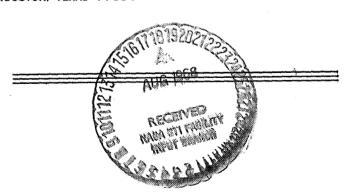


# MATHEMATICAL RESEARCH, INC.

HOUSTON, TEXAS 77004



NASA-(R-92216

METHODS OF THE GALERKIN TYPE

THIS REPORT WAS PREPARED WITH PARTIAL SUPPORT OF NASA-MSC(HOUSTON)
CONTRACT NAS 9-6964

# TABLE OF CONTENTS

CHAPTER		PAGE
I.	INTRODUCTION	1
II.	GENERAL APPROXIMATION THEORY	4
	A. General Definitions	14
	B. Theorems	13
	C. Summary	18
III.	CLASSIFICATION OF METHODS OF THE GALERKIN TYPE	20
	A. Definitions	21
	B. Specific Methods	23
	C. Summary	31
IV.	RESULTS ON METHODS OF THE GALERKIN TYPE	36
	A. Definitions	36
	B. Fundamental Theorems	38
	C. Improving Approximations	45
V.	A RELATED METHOD	50
VI.	APPLICATIONS AND EXAMPLES	56
VII.	THE NONLINEAR CASE	90
VIII.	APPLICATIONS AND ANALYSES OF METHODS OF THE	
	GALERKIN TYPE AVAILABLE IN THE LITERATURE	94
	A. Ordinary Differential Equations	94
	B. Partial Differential Equations	99
	C. Method Studies	102

#### Chapter I

#### INTRODUCTION

In his paper published in 1915, B. G. Galerkin [46] solved a number of problems on the equilibrium and stability of rods and plates by a method which has become associated with his name, the Galerkin method. Following this publication, a large amount of work has appeared in which the method has been extensively applied to the practical solution of diverse problems. The lack of proofs on convergence has never hindered its application. The first substantiation of the Galerkin method was given in 1940 by Repman [55] in application to Fredholm type integral equations, and in the same year an analogous result for a certain special ordinary fourth-order differential equation was obtained by Petrov [52]. In 1950 Mikhlin [51] obtained a fairly general sufficient condition for the convergence of the Galerkin method and demonstrated the application of this criterion to a number of problems which included the previously published results.

Through the years many modifications to the so-called classical Galerkin method have been suggested and used, in the vast majority of cases without proof of convergence. Justification is usually made on the basis of comparison of successive approximation or by comparison with experiment or other approximate solutions. It is only in recent years (since about 1957) that mathematicians (almost exclusively Russian) have begun to use functional analysis to obtain proofs for approximation methods and, in particular, the Galerkin and related methods.

The classical Galerkin method can briefly be described as follows:

Let A be a linear operator defined on a set  $D_{\!\!\!A}$  which is dense in some separable Hilbert space where a solution to the equation

$$Au - f = 0 \tag{1.1}$$

is required. A set of linearly independent elements  $\{\phi_n\}$  in  $D_A$  are selected and are called coordinate functions. For each integer n, an approximate solution of (1.1) is constructed in the form of a linear combination of the first n coordinate functions with constant coefficients, that is,

$$u_n(X) = \sum_{k=1}^n a_k \phi_n(X)$$
 (1.2)

The coefficients  $a_n$  are determined by replacing u by  $u_n$  in (1.1) and requiring the left-hand side of (1.1) to be orthogonal to the first n coordinate functions  $\phi_1, \ldots, \phi_n$ . This leads to the following system of n linear equations in the n unknowns  $a_k$ :

$$\sum_{k=1}^{n} (A \phi_{k}, \phi_{j}) a_{k} = (f, \phi_{j}) \qquad j = 1, ..., n$$
 (1.3)

Closely connected with the Galerkin method is the method of weighted residuals. In this method, a second set of linearly independent functions  $\{\psi_n\}$  in  $D_A$ , called weight functions, is selected. The  $n^{\text{th}}$  approximate solution is again given by (1.2), but the coefficients in this case are obtained as the solution of the following system: .

$$\sum_{k=1}^{n} (A \varphi_k, \psi_j) a_k = (f, \psi_j) \qquad j = 1, ..., n$$
 (1.4)

In connection with these and similar methods there are several outstanding and very practical problems. It is desireable to have some definition of an optimal approximation to the exact solution using methods of this type. Once such a definition is given, the question of selecting a set of optimal coordinate functions and a set of optimal weight functions remains. In addition, the question of the convergence of the approximate solutions to the exact solution is ever present. Furthermore, from a practical point of view, it is important to develop computational techniques which make efficient use of these methods as well as to anticipate any computational problems involved.

The purpose here is to attack the above problem areas as well as to present a concise, lucid guide to efficient and economical application of these methods and their modifications. In addition a rather complete classification of Galerkin type methods is given with existing convergence theorems. Finally, the last chapter includes a bibliography which represents a rather extensive search. Most of the references not included are those which make a straight-forward application of a method.

#### Chapter II

# GENERAL APPROXIMATION THEORY

It is the purpose of this chapter to state the general definitions used in the sequel and to present an outline of approximation theory [49] from a somewhat general point of view which will include some later methods as special cases.

#### A. General Definitions

A <u>linear space</u> over the field of real (complex) numbers is a set

V of elements called points or vectors satisfying the following axioms.

- (1) To every pair, f and g, of elements in V there corresponds an element f + g called the sum of f and g. With respect to the operation +, V is an abelian group, i.e.
  - (a) + is commutative, f + g = g + f,
  - (b) + is associative, f + (g + h) = (f + g) + h,
- (c) there exists in V a unique element O such that f + O = f for every f in V, and
- (d) for each f in V there corresponds a unique element -f in V such that f + (-f) = 0.
- (2) To every pair,  $\alpha$  and f, where  $\alpha$  is a real (complex) number and f is an element of V, there corresponds an element  $\alpha$ f in V, called the product of  $\alpha$  and f, such that

- (a) scalar multiplication is associative,  $(\alpha\beta)f = \alpha(\beta f)$ ,
- (b) lf = f, Of = O for every f in V,
- (c) scalar multiplication is distributive with respect to vector addition,  $\alpha(f+g) = \alpha f + \alpha g$ , and
- (d) multiplication by vectors is distributive with respect to scalar addition,  $(\alpha + \beta)f = \alpha f + \beta f$ .

Elements  $f_1, \ldots, f_n$  in V are <u>linearly independent</u> if the relation

$$\alpha_1 f_1 + \dots + \alpha_n f_n = 0$$

holds only for  $\alpha_1 = \ldots = \alpha_n = 0$ ; otherwise  $f_1, \ldots, f_n$  are <u>linearly dependent</u>.  $\alpha_1 f_1 + \ldots + \alpha_n f_n$  is called a <u>linear combination</u> of the elements  $f_1, \ldots, f_n$ . A linear space V is <u>n-dimensional</u> if V contains n linearly independent elements and if any set of n+1 elements in linearly dependent. If for each positive integer n, V contains n linearly independent elements, then V is <u>infinite dimensional</u>. A non-empty subset M of V is a <u>subspace</u> if for every pair,  $f_1$  and  $f_2$ , of elements of M, every linear combination,  $\alpha_1 f_1 + \alpha_2 f_2$ , is also contained in M. If  $f_1, \ldots, f_n$  are linearly independent in V, the space of all linear combinations of  $f_1, \ldots, f_n$  is called the <u>subspace spanned</u> by  $f_1, \ldots, f_n$  and is denoted by  $S_n(f)$  and  $f_1, \ldots f_n$  is called a <u>basis</u> for  $S_n(f)$ .

A linear space V is an <u>inner product space</u> if for each pair of elements, f and g, in V there is a real (complex) number (f, g),

called the inner (or scalar) product of f and g which satisfies the following conditions:

- (1)  $(g, f) = (\overline{f, g})$ , the complex conjugate of (f, g),
- (2)  $(\alpha_1 f_1 + \alpha_2 f_2, g) = \alpha_1 (f_1, g) + \alpha_2 (f_2, g),$
- (3)  $(f, f) \ge 0$ , and
- (4) (f, f) = 0 if and only if f = 0.

The norm of the element f in V, denoted by I f I, is defined by

$$\|f\| = \sqrt{(f, f)}$$

The following properties of the norm can be derived from the properties of the inner product.

- (1)  $\|\alpha f\| = \|\alpha\| \|f\|$
- (2)  $|(f, g)| \le ||f|| ||g||$ , the Cauchy inequality
- (3)  $\|f + g\| \le \|f\| + \|g\|$ , the triangle inequality

Two elements f and g are said to be orthogonal if

$$(f, g) = 0.$$

An element f is normalized if

$$||f|| = 1.$$

An inner product space  $\,V\,$  becomes a metric space if the distance between two elements, f and g, in  $\,V\,$ , denoted by  $\,d(f,\,g)$  is defined as

$$d(f, g) = ||f - g||.$$

From the properties of the norm, it can be shown that the function d satisfies the usual conditions for a metric, i.e.

(1) 
$$d(f, g) = d(g, f)$$

(2) 
$$d(f, g) \le d(f, h) + d(h, g)$$

(3) 
$$d(f, g) > 0$$

(4) 
$$d(f, g) = 0$$
 if and only if  $f = g$ .

If V is a metric space and  $\left\{f_n\right\}_{n=1}^{\infty}$  is a sequence of elements of V, then f in V is said to be the <u>limit of the sequence</u>  $f_n$ , written  $f_n \to f$  or  $\lim_{n \to \infty} f_n = f$ , if

$$\lim_{n\to\infty}d(f_n, f)=0.$$

The sequence  $\left\{f_n\right\}_{n=1}^{\infty}$  is said to <u>converge</u> to f. This type of convergence is called <u>strong convergence</u> or <u>norm convergence</u>. The sequence  $\left\{f_n\right\}_{n=1}^{\infty}$  is said to be a <u>Cauchy (or fundamental)</u> sequence, if for each  $\epsilon > 0$  there is an integer  $N_{\epsilon}$  such that for all m,  $n > N_{\epsilon}$ 

$$d(f_m, f_n) \leq \epsilon.$$

A metric space V is called <u>complete</u> if every Cauchy sequence converges to an element of V.

A <u>Hilbert space</u> H is an infinite dimensional inner product space which is a complete metric space with respect to the metric induced by the inner product. H is a <u>separable</u> Hilbert space if H contains a sequence of elements  $\left\{f_n\right\}_{n=1}^{\infty}$  such that  $\overline{S(f)}$ , the <u>closure</u> of S(f), that is, the subspace spanned by the set  $\left\{f_n\right\}$  along with limits of all Cauchy sequences in S(f), is the space H. Thus if H is a separable Hilbert space, each element f of H can be approximated with arbitrarily prescribed accuracy by a finite linear combination of the elements  $\left\{f_n\right\}_{n=1}^{\infty}$ .

The space  $L_2[a,b]$  of functions which are square integrable (in the Lebesque sense) on [a,b] is an example of a separable Hilbert space. In this case, denoting by  $\mu$  the measure on [a,b], the inner product of f and g in  $L_2[a,b]$  is defined as

$$(f, g) = \int_{a}^{b} fg d\mu$$

and the norm of f thus becomes

$$\| f \| = \sqrt{(f, f)} = \left[ \int_a^b f^2 d\mu \right]^{\frac{1}{2}}.$$

Norm convergence in L2 is called mean convergence.

Let G be a subspace of a separable Hilbert space H. It can be shown that to each element h in H there corresponds a unique element

g in G such that

$$\|h - g\| = \inf_{g' \in G} \|h - g'\|.$$

Furthermore, h - g is orthogonal to the subspace G, i.e., for each g' in G

$$(h - g, g^{\dagger}) = 0,$$

and h can be represented in the form

$$h = g + f$$

where  $g \in G$  and f is orthogonal to G. The element g is called the orthogonal projection of h on the subspace G.

A set 
$$M = \{f_n\}_{n=1}^{\infty}$$
 which satisfies

(1) 
$$(f_m, f_n) = \begin{cases} 0 \text{ for } m \neq n \\ 1 \text{ for } m = n, \text{ and } \end{cases}$$

(2) 
$$(f_n, f_n) = 1 \text{ for } n = 1, 2, \dots$$

is called an <u>orthonormal system or sequence</u>. From a given set of linearly independent (finite or countable) set of elements,  $\{g_n\}$ , it is possible to construct an orthonormal set,  $\{f_n\}$ , so that each  $f_n$  is a linear combination of the  $g_m$ ,  $m=1,\ldots,n$ . One such method of orthonormalization is the Gram-Schmidt process. An orthonormal (linearly independent) system is said <u>complete</u> if it is not contained in any larger orthonormal (linearly independent) set.

Let D denote a subspace of the separable Hilbert space H. A function T which relates to each element f in D a particular element Tf = g in H is called an operator in the space H with domain D. The set

$$\Delta_{T} = \{g \in H : Tf = g, f \in D\}$$

is called the <u>range</u> of T. The domain of T is denoted by  $D_T$ . The operator which maps each element onto itself is called the <u>identity</u> operator and is denoted by E.

If the operator T is one-to-one, that is, for each distinct pair,  $f_1 \text{ and } f_2, \text{ in } D_T, \text{ } Tf_1 \neq Tf_2 \text{ in } \Delta_T, \text{ then T has an } \underline{\text{inverse}} \text{ denoted}$  by  $T^{-1}$  which maps  $\Delta_T$  onto D. Furthermore,  $T^{-1}g = f$  if and only if Tf = g.

Let S and T be two operators such that  $\Delta_T \cap D_S$  is not empty. Then the <u>product of S and T</u> is defined as

$$STf = S(Tf).$$

An operator T is linear if for each pair f and g in  $\textbf{D}_{T}$  and scalars  $\alpha$  and  $\beta$ 

$$T(\alpha f + \beta g) = \alpha T f + \beta T g.$$

The norm of an operator T, denoted by ||T||, is defined as

as f ranges over  $D_T$  and  $\|f\| \le 1$ .

An operator T is bounded if ||T|| is finite.

The operator T is said to be continuous at a point  $f_0$  in  $D_T$  if for each  $\varepsilon>0$  there exists a  $\delta=\delta(\varepsilon)>0$  such that for each f in  $D_T$  with

$$\parallel f$$
 -  $f_{_{\scriptsize O}} \parallel$  <  $\delta$  we have  $\parallel Tf$  -  $Tf_{_{\scriptsize O}} \parallel$  <  $\varepsilon.$ 

If T is an arbitrary bounded linear operator on H, there exists a unique bounded linear operator  $T^*$  such that  $\|T\| = \|T^*\|$  and for each f and g in H

$$(f, T_g) = (T*f, g).$$

The operator  $T^*$  is called the <u>adjoint</u> of T. It is easy to see that  $(T^*)^* = T$ . If  $T^* = T$ , then T is said to be <u>self-adjoint</u>.

(Tf, f) 
$$\geq 0$$
.

If (Tf, f) = 0 implies f = 0, then T is said to be <u>positive definite</u>. The <u>energy product</u> of two elements f and g in  $D_T$  is (Tf, g), and the <u>energy norm</u>, denoted by  $\|f\|_T$  is

$$\|f\|_{T} = \sqrt{(Tf, f)}.$$

A sequence of elements  $\left\{f_n\right\}_{n=1}^{\infty}$  in  $D_T$  is said to converge in energy to the element f in  $D_T$  if

$$\lim_{n\to\infty} \|f_n - f\|_T = 0.$$

It can be shown that energy convergence implies mean convergence [38].

A linear operator T is said to be symmetric if  $D_{\rm T}$  is dense in H and for f and g in  $D_{\rm m}$ 

$$(Tf, g) = (f, Tg).$$

It is easily verified that a symmetric operator T such that  $\mathbf{D}_{\mathrm{T}} = \mathbf{H}$  is self-adjoint.

A set C is <u>compact</u> if every infinite sequence in C contains a convergent subsequence. A linear operator T such that  $D_T = H$  is said to be <u>completely continuous</u> if it maps each bounded set into a set which is compact.

Let H be a separable Hilbert space and G be a subspace of H such that each element h in H has a unique representation in the form

$$h = g + f$$

where g is the orthogonal projection of h on G. The operator which maps h into its orthogonal projection g on G is called the <u>orthogonal projection operator</u> and is denoted by  $P_G$  or simply P if the subspace G has been previously specified. An orthogonal projection operator is linear, bounded, and ||P|| = 1. In addition,  $P^2 = P$  and  $P^* = P$ .

A complex number  $\lambda$  is called an <u>eigenvalue</u> of the linear operator T if there exists a element  $f \neq 0$  such that

$$Tf = \lambda f$$
.

The vector f is called an <u>eigenvector</u> belonging to the eigenvalue  $\lambda$ .

A subspace G of H is called an <u>invariant subspace</u> of the operator T if every element of  $D_{\underline{T}}$  in G is mapped by the operator T into an element also in G.

Given a complete orthonormal set  $\{f_n\}$  in  $\overline{S(f)}$  each function g has the unique representation

$$g = \sum_{i=1}^{\infty} (g, f_i) f_i$$

and this is called the <u>Fourier series</u> of g relative to the set  $\{f_n\}$ . The orthogonal projection of g on  $S_n(f)$  can be written

$$\sum_{i=1}^{n} (g, f_i) f_i.$$

#### B. Theorems

Let  $\overline{X}$  be a complete subspace of the normed space X and P denote a projection from X onto  $\overline{X}$ , i. e.

$$P(X) = \overline{X}; P^2 = P.$$

We first consider an equation of the form

$$Kx = x - \lambda Hx = y \tag{2.1}$$

in the space X, called the exact equation, and an equation

$$\overline{K}\overline{x} \equiv \overline{x} - \lambda \overline{H}\overline{x} = Py$$
 (2.2)

in  $\overline{X}$ , called the <u>approximate</u> <u>equation</u>. The operators H and  $\overline{H}$  are linear.

In order to make the rather abstract setting more easily understood, we will take an example to illustrate the theory as we proceed.

Let  $\{\phi_i\}$  be an orthonormal sequence in  $L_2$ . Let  $X=\overline{S(\phi)}$  and  $\overline{X}=S_n(\phi)$ . Then the operator P could be defined as the orthogonal projection from X onto  $\overline{X}$ .

The following conditions will play an important role in the sequel.

- II. For every x in X there is an  $\overline{x}$  in  $\overline{X}$  such that  $\| \operatorname{Hx} \overline{x} \| \leq \eta_1 \, \| \, x \, \| \, .$
- III. For each y in X there is a  $\overline{y}$  in  $\overline{X}$  such that  $\|y-\overline{y}\| \leq \eta_2 \|y\|$ .
- IV. If a solution exists for equation (2.2) then the solution is unique.

In our example, condition I could be satisfied if  $\overline{H}=PH$ . In such case  $\eta=0$  is an acceptable choice. Furthermore, conditions II and III would be satisfied provided n is sufficiently large since the sequence  $\{\phi_i\}$  is complete.

Condition IV will always hold for completely continuous operators.

We hope to be able to solve the approximate equation and have an approximation to the exact equation. Suppose that we have a sequence of approximate equations and solutions obtained from these. Then the space

 $\overline{X}$ , together with the operators  $\overline{H}$  and  $\overline{P}$  and the constants in conditions I, II, and III depend on an index n, which we shall omit for the sake of simplicity.

The obvious question is whether or not such a sequence of approximations converges to the exact solution.

# Theorem 2.1. If the following conditions are satisfied:

- (1) K has a linear inverse
- (2)  $\overline{K}$  satisfies IV for each n
- (3) I, II, and III hold for each n
- (4)  $\lim_{n\to\infty} \eta = 0$ ;  $\lim_{n\to\infty} \eta \|P\| = 0$ ;  $\lim_{n\to\infty} \eta_2 \|P\| = 0$

then the approximate equations are soluable for sufficiently large n and

$$\lim_{n\to\infty} \|x - \overline{x}_n\| = 0.$$

These sufficient conditions for convergence can be phrased in another form when  $\overline{H}=PH$ .

# Theorem 2.2. If $\overline{H} = PH$ and the following conditions are satisfied:

- (1) X is complete
- (2)  $\lim_{n\to\infty} P_n x = x$
- (3) H is completely continuous
- (4)  $\lambda$  is not the reciprocal of an eigenvalue of H

then the approximate solutions converge to the exact solution.

Corollary: The characteristic values of H are limits of sequences of characteristic values of  $\overline{H}_n$ .

In our example, X is complete and the orthogonal projectors converge to the identity so (1) and (2) of Theorem 2.2 are satisfied.

We shall now consider equations in which the left-hand side does not contain the identity operation explicitly and moreover represents an operation from the original space not into itself but into another normed space. In spite of these differences as compared with the equations discussed above, these are reducible to the former type.

We suppose that X and Y are normed spaces, from each of which complete subspaces  $\overline{X}$  and  $\overline{Y}$ , respectively, have been chosen. We suppose that F is a linear operator projecting Y onto  $\overline{Y}$ .

As before, we consider two equations, the exact

$$K_1 x = Gx - \lambda Tx = y_1$$
 (2.3)

and the corresponding approximate equation

$$\overline{K}_1 \overline{x} \equiv G \overline{x} - \lambda \overline{T} \overline{x} = F y_1.$$
 (2.4)

Here, G and T (and  $K_1$ ) are linear operators mapping X into Y while  $\overline{T}$  (and  $\overline{K}_1$ ) are linear operators from  $\overline{X}$  into  $\overline{Y}$ . We further assume that G has a linear inverse and establishes a one-to-one correspondence between  $\overline{X}$  and  $\overline{Y}$ . In this setting conditions I through III above become:

Ib. For every x in X  $\| FT\overline{x} - \overline{T}\overline{x} \| < \mu \| | \overline{x} | \|$ 

IIb. For every x in X there is a  $\overline{y}$  in  $\overline{Y}$  such that  $\|Tx - \overline{y}\| \le \mu_{\gamma} \|x\|$ 

IIIb. For each  $y_1$  in Y there is an element  $\overline{y}_1$  in  $\overline{Y}$  such that  $\|y_1-\overline{y}_1\|\leq \mu_2$   $\|y_1\|$ 

Writing, under the above assumptions, (2.1) and (2.2) as

$$Kx \equiv G^{-1}K_1x \equiv x - \lambda G^{-1}Tx = G^{-1}y_1$$
 (2.5)

$$\overline{K}\overline{x} \equiv G_1^{-1}\overline{K}_1\overline{x} \equiv \overline{x} - \lambda G^{-1}\overline{T}\overline{x} = G^{-1}Fy_1$$
 (2.6)

it can be shown that conditions I through III hold.

The preceding theorems can be restated in this setting.

### Theorem 2.3. If the following conditions are satisfied:

- (1) K<sub>1</sub> has a linear inverse,
- (2)  $\overline{K}_1$  satisfies IV for each n,
- (3) Ib, IIb, IIIb are satisfied for each n,
- (4)  $\lim_{\mu \to \infty} \mu = 0$ ;  $\lim_{\mu \to \infty} \mu_1 \| F \| = 0$ ;  $\lim_{\mu \to \infty} \mu_2 \| F \| = 0$ ,

then the approximate equations are soluable for sufficiently large n, and the sequence of approximate solutions converges to the exact solution.

As before when the approximate equation is obtained by projection of the exact equation, the theorem can be formulated differently.

Theorem 2.4. If  $\overline{T} = FT$  and the following conditions are satisfied:

- (1) Y is complete
- (2)  $\lim_{n \to \infty} F_n x = x$
- (3)  $G^{-1}T$  is completely continuous
- (4)  $K_1^{-1}$  exists

then the approximate equations are soluable for sufficiently large n, and the approximate solutions converge to the exact solution.

#### C. Summary

From the above we conclude that when presented an operator equation to solve, we must decide into which general category it falls. That is, decide whether it is of type 1 with the identity explicit or type 2. Then several other definitions must be made. There is a great deal of flexibility in these, but certain choices are sometimes more judicious than others. There is no general rule since the choice is influenced by the character of the operators as well as the boundary conditions and the particular method of solution employed.

The following is an outline of the situation.

Type 1.  $x - \lambda Hx = y$ 

Define: (1) the space X and the subspace  $\overline{X}_n$ 

- (2) the projection P<sub>n</sub>
- (3) the operator  $\overline{H}_n$

At this point the approximate equation is determined, and it may be possible to check conditions I through IV. It must then be determined if either Theorem 2.1 or Theorem 2.2 applies to prove that the approximate solutions converge.

Regarding type 2, it should be pointed out that the first step here is to write the given equation in the proper form. There may be more than one way to do this, and one version may be more amenable than another. Again, there is no general principle. However, it must be kept in mind that G must be one-to-one and have a linear inverse.

Type 2.  $Gx - \lambda Tx = y$ 

Define: (1) the spaces X,  $\overline{X}$ , Y and  $\overline{Y}$ 

- (2) the projection F
- (3) the operator  $\overline{T}$

As above the approximate equation is now determined and, perhaps, condition Ib through IIIb, and IV can be checked. Also it remains to determine if the theorems apply to prove convergence.

The above procedures are exhibited in the examples.

#### Chapter III

#### CLASSIFICATION OF METHODS OF THE GALERKIN TYPE

Since Galerkin first used his technique to solve a problem in mechanics, a vast number of modifications have been suggested. The purpose of this section is to classify these suggested modifications and to outline their application. Certain of these fall into the general approximation theory considered in the last chapter. Some of them have been rather widely used but are without convergence proofs. These will be carefully noted along with the limitations of the other methods.

We first state the definitions applicable to this chapter, then classify and describe the methods and in the final section summarize the results and draw certain conclusions.

The methods all involve the solution of the equation

$$L(u) = 0 \tag{3.1}$$

where u is defined on a domain D bounded by C, L is an operator in a separable Hilbert space and on C

$$B_{i}(u) = g_{i}, i = 1, ..., p$$
 (3.2)

are the appropriate number of boundary conditions. It should be noted here, and will be pointed out below, that the methods discussed are primarily designed for boundary valve problems and certain careful modifications must be made to even apply any of these methods to initial value or other types of problems. The methods to be described involve the selection

of two linearly independent sets of functions: the first called coordinate functions  $\left\{v_i\right\}_{i=1}^{\infty}$  and the second called weight functions  $\left\{w_i\right\}_{i=1}^{\infty}$ .

An n approximation to the solution is obtained in the form

$$u_n = f(c_1, ..., c_n, v_1, ...v_n)$$
 (3.3)

where f is a function of some selected form (e.g., a linear combination) and the functions (perhaps constants) c are determined by one of the procedures described below.

#### A. Definitions

<u>Definition 3.1:</u> The <u>operator residual</u>, denoted by  $R_L(u_n)$ , is the result of operating by L on  $u_n$ . Thus  $R_L(u_n) = L(u_n)$ .

Definition 3.2: The component boundary residuals, denoted by  $R_{B_i}(u_n)$ , are the result of substituting  $u_n$  for u in the boundary conditions. Thus  $R_{B_i}(u_n) = B_i(u_n) - g_i$ ,  $i = 1, \ldots, p$ .

Definition 3.3: The boundary residual, denoted by  $R_B(u_n)$ , is the sum of the absolute values of the component boundary residuals. Thus  $R_B(u_n) = \sum_{i=1}^p \left| R_{B_i}(u_n) \right|$ .

Definition 3.4: A method of solution is called an interior method provided that the coordinate functions are chosen to satisfy the boundary conditions and the functions  $\{c_i\}$  are determined by requiring that the operator residual be orthogonal to the first n weight functions throughout the region D. That is

$$(R_L(u_n), w_j) = 0$$
  $j = 1, ..., n.$  (3.4)

Definition 3.5: A method of solution is called a <u>boundary method</u> provided that the coordinate functions are chosen to satisfy the operator equation, and the functions  $\{c_i\}$  are determined by requiring that the boundary residual be orthogonal to the first n weight functions on the boundary C. That is

$$(R_B(u_n), w_j) = 0 \quad j = 1, ..., n.$$
 (3.5)

<u>Definition 3.6</u>: A method of solution is called a <u>mixed</u> <u>method</u> if it is neither interior nor boundary.

It should be noted that generally in the mixed method described above both the boundary and operator residuals must be considered. If we require that both the boundary residual and the operator residual be orthogonal to the first n weight functions, we will obtain 2n equations for the n unknowns  $\{c_i\}$ . It has been suggested that such a situation be remedied by:

- (1) Discarding some of the equations; the choice is arbitrary as long as the equations are independent although consideration is given to those characteristics of the problem which are most important [29]. Exactly how such a decision should be made is not stated.
- (2) Using  $\frac{n}{2}$ , instead of n, weight functions (n must, of course, be even) [13].
- (3) Adding the residuals and making this sum orthogonal to the first n weight functions [36].

None of these suggestions are proved to be valid and are merely ways of obtaining n equations in n unknowns. The sole motivation for the rather arbitrary choice seems to be a matter of taste.

#### B. Specific Methods

There are four fundamental classifications of methods of the Galerkin type.

- I. <u>General</u>: This is simply to solve the system (3.4) using the form (3.3) as it stands. Unfortunately, it can be quite difficult to achieve such a solution, and a certain amount of ingenuity may be required to select an appropriate form for f, to satisfy the boundary conditions, etc.
- II. Shulesko Technique: This is a mixed method wherein p + 1 sets of coordinate functions are chosen,  $\left\{v_i\right\}_{i=1}^{\infty}$  and  $\left\{t_{ki}\right\}_{i=1}^{\infty}$ ,  $k=1,\ldots p$ . The n<sup>th</sup> approximation is sought in the form

$$u_n = \sum_{k=1}^{p} \sum_{i=1}^{n} a_{ki} t_{ki} + \sum_{i=1}^{n} b_i v_i$$

Furthermore two complete sets of weight functions  $\{s_j\}$ ,  $\{w_j\}$  are selected. The  $\{a_{ik}\}$  are then determined by making the component boundary residuals orthogonal to the first n weight functions  $s_j$ . The  $b_i$  are determined by making the operator residual orthogonal to the first n weight functions  $w_j$ . [41].

III. Method of Weighted Residuals: A set of functions  $\left\{v_i\right\}_{i=1}^{\infty}$  is chosen to satisfy the boundary conditions. The  $n^{th}$  approximation

takes the form  $u_n = \sum_{i=1}^n c_i v_i$  and the  $\{c_i\}$  are determined by making the operator residual orthogonal to the first n of any set of linearly independent weight functions  $\{w_j\}_{j=1}^{\infty}$  [48].

IV. Classical Galerkin: A set of coordinate functions of the required number of variables, linearly independent and complete, is chosen to satisfy the boundary conditions. The  $n^{th}$  approximation has the form  $u_n = \sum_{i=1}^n c_i v_i$  where the  $\{c_i\}$  are constants determined by making the operator residual orthogonal to the set of coordinate functions. Thus the coordinate functions themselves become the weight functions [38].

Notice that II includes III in the case in which  $\{t_{ki}\} \equiv \{v_i\}$  and satisfy the boundary conditions; it includes IV in the case in which  $\{t_{ki}\} \equiv \{v_i\}$  and satisfy the boundary conditions and  $\{w_j\} \equiv \{v_i\}$ . In addition, III includes IV when  $\{v_i\} \equiv \{w_j\}$ .

There are several slight modifications and variations on the above theme. These depend on the specific characteristics of the problem to be solved; for example, whether or not the operator is linear, the boundary conditions homogeneous, etc. Some of these are indicated below.

#### (1) Ames Method

When the problem is an initial value problem, where the range of t is the non-negative real numbers, the n<sup>th</sup> approximation is taken in the

form  $u_n = \sum_{i=1}^n c_i v_i$  where the coordinate functions satisfy the initial conditions. The operator residual is then made orthogonal to a set of weight functions (according to III or IV) on the interval [0, T] where T is regarded as the time when steady state is achieved. Hopefully, T may be estimated from the physical nature of the problem [31]. See Chapter V for a further discussion of this.

# (2) Bickley Method

When a linear problem involves the variables x and t, the  $n^{th}$  approximation may be taken in the form  $u_n = \sum_{i=1}^n c_i(t)v_i(x)$ . Thus the members of  $\{c_i(t)\}$  are unknown functions of t and may be determined from the conditions of III or IV [18].

# (3) Crandall Method

When L is a partial differential operator, it is possible to reduce the problem to an ordinary differential system if we are given in addition to (3.1) and (3.2) an appropriate number of initial conditions:

$$I_k(u) = h_k$$
  $k = 1, ..., r$  (3.6)

The n<sup>th</sup> approximation is taken in the form  $u_n = \sum_{l=1}^n c_i(t)v_i$  where  $\left\{v_i\right\}_{i=1}^{\infty}$  is linearly independent, complete, and satisfies the boundary conditions but not the initial conditions or the differential equation. Thus, in addition to the operator residual, we define an <u>initial residual</u> to be:

$$R_{\underline{I}}(u_n) = \sum_{k=1}^{r} \left| h_k - I_k \left( \sum_{i=1}^{n} e_i(0) v_i \right) \right|.$$

Making these residuals orthogonal to the first n of a complete set of weight functions  $\{w_j\}_{j=1}^{\infty}$ , we obtain the two systems of equations:

(a) 
$$\int_{\mathbf{m}} \mathbf{R}_{\mathbf{I}}(\mathbf{u}_{\mathbf{n}}) \mathbf{w}_{\mathbf{j}} d\mathbf{T} = 0$$

and

(b) 
$$\int_{D} R_{L}(u_{n})w_{j}dD = 0$$

for j = 1, ..., n. From equation (b) we get a system of ordinary differential equations in the  $\{c_i(t)\}$ , and equation (a) will specify initial conditions for them, allowing solution either exactly or approximately [34].

# (4) Generalized Moments Method

W. V. Petryshyn [53] has suggested a method for the solution of the eigenvalue problem

$$T_{ij} - \lambda Su = 0 (3.7)$$

where the operators T and S are linear, unbounded and nonsymmetric, and T is K-positive definite. An operator is said to be K-positive definite provided there is a closeable operator K with  $D_{K} \geq D_{T}$  mapping  $D_{T}$  onto a dense subset of H and two positive constants  $\alpha_{1}$  and  $\alpha_{2}$  such that

(Tu, Ku) 
$$\geq \alpha_1 \| u \|^2$$
  $u \in D_T$ 

and

|| Ku|| 
$$^2 \leq \alpha_2$$
 (Tu, Ku) ueD<sub>T</sub>

An operator K is closeable if  $\{u_n\}$  converging to zero and  $\{Ku_n\}$  converging to f imply f is zero. He also suggested the method for

$$Au + Bu = f (3.8)$$

where A is K-positive definite, and B is a linear unbounded operator [54]. The method is the method of weighted residuals where the weight functions are chosen as the K-images of the coordinate functions. The author points out that if T (or A) is self-adjoint, the method reduces to the classical Galerkin method (i.e., K = I).

Certain of the above are special cases of the general theory of the preceding chapter. The classical Galerkin method is a special case in the following way. If the operator L has the form of type 1, the space X is taken as the completion of the subspace of  $L_2$  spanned by coordinate functions and the subspace  $\overline{X}$  as the subspace spanned by the first n coordinate functions. Further, P is taken as the orthogonal projection operator from X onto  $\overline{X}$  and  $\overline{H}$  is defined as PH. The approximate equation then gives, on taking the inner product with  $\phi_1$ ,

$$(\overline{x} - \lambda \overline{H} \overline{x}, \varphi_j) = (Py, \varphi_j)$$
  $j = 1, \dots n.$ 

Simplifying, using the self-adjointness of P and the fact that  $P(\phi_j) = \phi_j, \ j \le n, \ \text{gives}$ 

$$(\overline{x}, \varphi_j) - \lambda (\overline{H} \overline{x}, \varphi_j) = (y, \varphi_j) \quad j = 1, \dots n$$
 (3.6)

which is exactly the Galerkin system of equations. Since

$$\overline{x} = \sum_{i=1}^{n} a_{i} \phi_{i}$$

equation (3.6) gives n equations in the n unknowns  $a_i$ .

If the operator is of type 2, then X and  $\overline{X}$  are taken as above. Further, Y and  $\overline{Y}$  are taken as G(X) and  $G(\overline{X})$ , respectively. The operator F is taken as the orthogonal projection from Y onto  $\overline{Y}$ . Then the approximate equation is

$$G\overline{X} - \lambda FT\overline{X} = Fy.$$

As above this gives rise to the system

(Gz, 
$$\phi_j$$
) -  $\lambda$ (FT $\lambda$ ,  $\phi_j$ ) = (Fy,  $\phi_j$ )

This system determines the nth approximation.

The method of weighted residuals can also be considered as a special case of the general theory under certain conditions. In this method the same approximate equations as above are employed, but the inner products are taken with weight functions  $\{w_i\}$  other than the coordinate functions. This gives

$$(\overline{x} - \lambda \overline{H} \overline{x}, w_j) = (Py, w_j)$$
  $j = 1, \dots n$ 

and

$$(G\overline{x}, w_j) - \lambda(FT\overline{x}, w_j) = (Fy, w_j)$$
 j = 1, ... n

for types 1 and 2, respectively. The self-adjointness of P (or F) can be invoked to write

$$(Py, w_{j}) = (y, Pw_{j})$$

However, unless  $w_j$  is in the span of the first n coordinate functions,  $P w_j \neq w_j$ . Therefore, it can be concluded that the method of weighted residuals is a special case of the general theory only when the span of the first n coordinate functions includes the first n weight functions. (See in this connection Theorem 4.5 of Chapter IV.)

This same problem arises when the coordinate functions are not orthogonal. In this case the projection used is not orthogonal and hence not self-adjoint. That is, the projection is defined as truncation of the expansion in terms of the coordinate functions at n terms. This means, in particular, that the right hand side of the Galerkin system should be  $(P\ y,\ \phi_j)$  and not  $(y,\ \phi_j)$  to be a special case of the general theory. It does, indeed, seem clear that a convergence proof using  $(P\ y,\ \phi_j)$  would imply convergence using the right hand side  $(y,\ \phi_j)$  since  $P\ y$  converges to y. However, the simple way out of this apparent dilemma about convergence as well as accuracy of corresponding approximations is to utilize the following theorem.

Theorem 3.1: If the coordinate functions  $\{\phi_i\}$  for the linear problem

$$Ay = f$$

are linearly independent but not orthogonal in  $L_2$  then the  $n^{th}$  approximation is the same as the  $n^{th}$  approximation using the sequence  $\left\{\phi_i^*\right\}$  obtained by orthonormalization of  $\left\{\phi_i\right\}$ .

<u>Proof</u>: Let G be the nonsingular  $n \times n$  matrix which takes  $\overline{\phi}$  to  $\overline{\phi}^*$ , i.e.,

$$\overline{\varphi}^* = \left\{ \begin{array}{c} \varphi_1^* \\ \vdots \\ \varphi_n^* \end{array} \right\} = G\overline{\varphi} = \left\{ \begin{array}{c} \varphi_1 \\ \vdots \\ \varphi_n \end{array} \right\}.$$

Let  $P = (a_i)$  be the  $n \times n$  matrix where  $a_i = (A\phi_i, \phi_j)$  and  $Q = (b_i)$  be the  $n \times n$  matrix where  $b_i = (A\phi_i^*, \phi_j^*)$ . In addition, let

$$y_n = \overline{c} \overline{\phi}$$
 and  $y_n^* = \overline{d} \overline{\phi}^* = \overline{d}G\overline{\phi}$ 

where  $\overline{c} = (c_1, \dots, c_n)$  and  $\overline{d} = (d_1, \dots, d_n)$ . The Galerkin systems are then

$$P\overline{e}' = R = \begin{cases} (f, \varphi_1), & Q\overline{d}' = S = \\ (f, \varphi_1), & Q\overline{d}' = S = \\ (f, \varphi_n), & (f_n, \varphi_n) \end{cases}$$

where prime denotes transpose.

If it can be shown that  $\overline{c} = \overline{d}G$ , the proof will be complete. We have

so that

$$L(\overline{\varphi}^*) \overline{d}^{*} = GL(\overline{\varphi})\overline{\varphi}^{*} = GL(\overline{\varphi})\overline{\varphi}^{*}G^{*}.$$

Integrating the extremes of this gives

$$Q = GPG^{\dagger}$$

so

$$P = G^{-1}QG^{-1}$$
.

Thus

$$e' = P^{-1}R = (G'Q^{-1}G) (G^{-1}S) = G'Q^{-1}S = G'\overline{d}'$$

and taking transposes gives the result.

The above theorem does answer the dilemma because orthonormal coordinate functions  $\{\phi_{{\tt i}}{}^*\!\}$  gives

$$(Pf, \phi_i^*) = (f, P\phi_i^*) = (f, \phi_i^*)$$

since P is an orthogonal projection, is self-adjoint and is, of course, the identity on its range. This means that the "correct" right hand side is used in the Galerkin system when the coordinate functions are orthogonal. However, by the theorem, the approximate solution with nonorthogonal linearly independent functions is the same even though the right hand side in this system is  $(f, \phi_i)$ .

#### C. Summary

The number of convergence proofs for the above methods is quite small. As noted above in certain cases a few techniques are covered by the earlier theorems from the general theory.

The motivation for accepting these methods probably stems from the following well known theorem.

Theorem 3.2. An orthogonal sequence of functions is complete if and only if the only function orthogonal to each member of the sequence is equivalent

to zero. This theorem implies that if

$$(L(u_n), \varphi_j) = 0$$

for all j then

$$L(u_n) = 0$$

i.e.,  $u_n$  is a solution, whenever  $\{\phi_j\}$  is a complete orthogonal sequence. Of course, only n values of j are used so that  $u_n$  is an "approximate" solution. In practice and in fact in most of the examples in the literature, the coordinate functions are neither orthogonal nor complete. Hence the appeal of this argument wanes.

The following is a summary of essentially all of the theorems proved in regard to convergence of the above methods.

Theorem 3.3. If the sequence of coordinate functions is complete, the classical Galerkin method converges in the mean to the exact solution of a Fredholm integral equation [38].

Theorem 3.4. Suppose Au = f has a unique solution and can be written as

$$A_{o}u + Ku = f$$

where  $A_0$  is symmetric and positive-bounded-below in a Hilbert space H. Then if  $A_0^{-1}$  K is completely continuous, the classical Galerkin method with a complete sequence of coordinate functions converges in energy to the exact solution. A similar result holds for the eigenvalue problem [38]. As an application of Theorem 3.4 the following result can be proven.

# Theorem 3.5. Consider the problem

$$(-1)^{m} u^{(2m)} - \lambda Ku = f$$

$$u(a) = u'(a) = \dots = u^{(m-1)} (a) = 0$$

$$u(b) = u'(b) = \dots = u^{(m-1)} (b) = 0$$

where K is a linear differential operator of order 2m-1 such that all coefficients are sufficiently differentiable. Suppose there is a unique solution. Then the classical Galerkin method with a complete sequence of coordinate functions converges to the exact solution. A similar result holds for the eigenvalue problem [38].

Theorem 3.6. Suppose, on some domain D,

$$-\sum_{i, k=1}^{m} \frac{\partial}{\partial x_{i}} \left( A_{i_{k}} \frac{\partial u}{\partial x_{k}} \right) + \sum_{i=1}^{m} B_{i} \frac{\partial u}{\partial x_{i}} + Cu = f$$

where u vanishes on the boundary of D, the coefficients may be variable, and the equation is elliptic. Then the classical Galerkin method with a complete sequence of coordinate functions converges in energy to the exact solution [38].

Theorem 3.7. Suppose, on the interval  $0 \le x \le \pi$  and for  $t \ge 0$ ,

$$\frac{\partial^2 u}{\partial x^2} - \frac{\partial u}{\partial t} - gu = f$$

where

$$u(x, 0) = u(0, t) = u(\pi, t) = 0,$$

f, g and their first two derivatives are continuous and  $f(0, 0) = f(\pi, 0) = 0$ . Then Bickley's variation of the classical Galerkin method

with {sin nx} as coordinate functions converges to the exact solution [47].

Theorem 3.7 can be slightly generalized. For example a first derivative term in x could be present and the coefficient of the first derivative in t could be a positive function, and the same conclusion will hold.

The following is a generalization of Theorem 3.4.

## Theorem 3.8. If the equation

Au + Bu = f

has a unique solution and  $T = A^{-1}B$  can be extended to a completely continuous operator then for sufficiently large n, there is a unique n<sup>th</sup> approximate solution from the generalized moments methods, and this sequence converges to the solution [54].

Recall that in using the generalized moments method A should be K-positive definite. The practical difficulty in making use of Theorem 3.8 is that there is no method to find K or even to decide if a given operator is K-positive definite.

Considering the above theorems, it is apparent that the convergence of the many suggested variations remains for the most part unproved. Papers dealing with a method per se are few compared with those that simply apply a method without regard to convergence. The acceptance of these techniques appears to stem from the fact that the solutions

"agree" in some sense with either experiment or approximate solutions obtained in other ways.

### Chapter IV

## RESULTS ON METHODS OF THE GALERKIN TYPE

Most of the theorems, definitions and computational techniques stated in this chapter could not be found in the literature. The theorems either generalize existing results or give new information. The techniques of solution of a problem are either a consequence of a theorem or the result of experience with examples.

We will consider throughout this chapter the problem

$$Ay = f (4.1)$$

subject to homogeneous boundary conditions where A is a self-adjoint linear operator defined in a separable Hilbert space. If the boundary conditions are not homogeneous a slight modification to the problem will usually make them so. In particular we investigate the question of selecting "optimal" sets of coordinate functions and weight functions for the method of weighted residuals.

Throughout this chapter, it is assumed that (4.1) has a unique solution which will be denoted by y\*. In particular, then, the operator L is one-to-one and has the property that for any linearly independent set  $\left\{\phi_i\right\}_{i=1}^n \quad \text{in } \ D_L, \text{ the set } \left\{L(\phi_i)\right\}_{i=1}^n \quad \text{is linearly independent in the range of L.}$ 

#### A. Definitions

operator A provided that the Galerkin method using this set of functions as coordinate functions gives the best  $L_2$  (energy) approximation in  $S_n(\phi)$  to y\*.

Definition 4.2. Suppose each of  $\{\phi_i\}_{i=1}^{\infty}$  and  $\{w_i\}_{i=1}^{\infty}$  is a linearly independent set of functions. Then  $(\{\phi_i\}, \{w_i\})$  is called an  $\underline{L}_2$  (energy) optimal weighted residual pair with respect to A provided this pair used in the method of weighted residuals gives the best  $\underline{L}_2$  (energy) approximation in  $\underline{S}_n(\phi)$  to  $\underline{y}$ .

Definition 4.3. The operators  $A_1$  and  $A_2$  with the same domain are strongly similar provided there is a linear operator P such that  $PA_1 = A_2P$ .

Definition 4.4. The operators  $A_1$  and  $A_2$  with the same domain are similar provided there exist continuous linear operators  $P_1$  and  $P_2$  such that  $A_1P_1 = A_2$  and  $A_2P_2 = A_1$ .

Definition 4.5. A pair of linearly independent sets of functions  $\{\{\phi_i\}, \{w_i\}\}$  is an  $\underline{L}_2$  (energy) almost optimal weighted residual pair with respect to  $\underline{A}$  provided there is a number  $\underline{C}$  such that for all i

$$\|A(w_i) - \varphi_i\| < C$$

where the norm used in the  $L_2$  (energy) norm. It is obvious that optimal implies almost optimal.

Definition 4.6. Suppose  $\{\phi_i\}$  is a linearly independent set of coordinate functions for the Galerkin method and  $y_n$  is the  $n^{th}$ 

approximate solution to (4.1). Let  $B_{\phi}$  denote the operator which maps the right hand side, f, to  $y_n$  and  $R_{\phi} = A^{-1} - B_{\phi}$ . Then  $\left\{\phi_i\right\}$  will be called an R-optimal Galerkin coordinate set if for all sets  $\left\{\psi_i\right\}$ 

$$\|R_{\varphi}\| \le \|R_{\psi}\|$$
.

The set will be called <u>almost R-optimal</u> if there is a positive number C such that

$$\|R_{\phi}\| \le C \inf \|R_{\psi}\|$$

where the inf is taken over all  $\left\{\psi_{\mathbf{i}}\right\}$  in the space.

In connection with the above definition it should be pointed out that examination of the norm of the remainder operator could very well not be an indication of the smallness of  $||y^* - y_n||$  since the operator norm is defined as

$$\| R_{\phi} \| = \inf \| R_{\phi} g \|$$

and there could be some function g such that

$$\|R_{\phi}\| < \|R_{\phi}g\| < \|R_{\phi}f\| = \|y^* - y_n\|.$$

### B. Fundamental Theorems

The only previously published attempt to find an optimal use of the methods considered herein is found in [45]. The approach used in this reference is to introduce the nonstandard energy norm. (See Definitions,) This allows certain interesting results to be established.

In [45] the problem (4.1) with homogeneous boundary conditions where A is a positive definite self-adjoint operator is investigated. They consider the domain of A to be  $L_2$ ,  $H_1$ , and  $H_2$  respectively, which are determined by the three inner products (u, v), (Au, v) and (Au, Av).

In this setting the following theorems are obtained by these authors.

Theorem 4.1. The best n<sup>th</sup> approximation in energy to the solution of (4.1) is given by the classical Galerkin method regardless of the coordinate system. If the coordinate functions are complete, then the approximations converge in energy to the solution.

Theorem 4.2. The set  $\{\phi_i\}$  is an R-optimal coordinate set for (4.1) if the span of the first n elements is the same as the span of the first n eigenfunctions of the problem

 $Au = \lambda u$ 

with homogeneous boundary conditions.

Theorem 4.3. An almost R-optimal set for a problem defined by an operator  $A_1$  is an almost R-optimal set for a problem defined by an operator  $A_2$  if  $A_1$  and  $A_2$  are similar.

These same theorems are also proved in the space obtained using the inner product (Au, Av).

It is clear that the optimal approximation to the solution y\* by a function in  $S_n(\phi)$  will be the orthogonal projection,  $y_n$ , of y\* onto  $S_n(\phi)$ . Also

$$y^* = y_n + y_n'$$
,  $(y_n', \phi_i) = 0$ ,  $i = 1$ , ... n.

in whatever space is being considered.

The following result from this investigation is in the spirit of the above theorems.

Theorem 4.4. In the norm induced by (Au, v) the pair ( $\{\phi_i\}$ ,  $\{A\phi_i\}$ ) is an optimal weighted residual pair for (4.1) for any linearly independent coordinate set  $\{\phi_i\}$ .

Proof: We know

$$[y_0, \phi_1] = [y_n, \phi_1].$$

That is,

$$(Ay_0, A\phi_1) = (Ay_n, A\phi_1).$$

However  $Ay_0 = f$  and letting  $A\phi_i = w_i$  gives

$$(f, w_i) = (Ay_n, w_i)$$

which is the weighted residual system. Notice that the self-adjointness of A is not used. A somewhat stronger and much more useful result than Theorem 4.4 is the following for self-adjoint operators.

Theorem 4.5. The coordinate set  $\{\phi_i\}$  is an  $L_2$  optimal Galerkin set for A in (4.1) if and only if  $S_n(\phi)$  is an invariant subspace of A.

Proof: Suppose  $S_n(\phi)$  is invariant under A. Then  $A\phi_j$  is in  $S_n(\phi),\ j\leq n$  and

$$(y_0, A\phi_j) = (y_n, A\phi_j).$$

Since A is self adjoint this gives

$$(Ay_0, \varphi_1) = (Ay_n, \varphi_1)$$

or

$$(f, \varphi_j) = (Ay_n, \varphi_j)$$

which says  $y_n$  is the  $n^{th}$  Galerkin approximation.

Suppose  $\textbf{S}_n(\phi)$  is not invariant under A. Then for some  $j \leq n,$   $\textbf{A}\phi_j$  is not in  $\textbf{S}_n(\phi).$  Then

$$(Ay_n, \phi_j) = (y_n, A\phi_j) \neq (y_o, A\phi_j) = (Ay_o, \phi_j) = (f, \phi_j).$$

so that  $y_n$  is not the  $n^{th}$  Galerkin approximation.

Corollary: The eigenfunctions of

$$Au = \lambda u$$

subject to homogeneous boundary conditions form an L<sub>2</sub> optimal Galerkin coordinate set.

Proof:  $S_n(\varphi)$  is invariant for each n and hence Theorem 4.5 gives the result.

Theorem 4.5 is useful since if a set of coordinate functions have been chosen, it can easily be determined if  $S_n(\phi)$  is an invariant subspace and hence if the Galerkin method yields the best  $L_p$ -approximation.

Once a set of coordinate functions is selected the question arises, "Is it possible to choose weight functions such that the pair of sequences will form an  $L_2$ -optimal weighted residual pair?" The next theorem answers this question in the affirmative under most circumstances. That is, for the problem (4.1) to have a solution, the function f, the right-hand side, must satisfy certain conditions. We will call such functions admissable. Then we have another theorem.

Theorem 4.6. Corresponding to any set  $\{\phi_i\}_{i=1}^{\infty}$  of linearly independent admissable functions which satisfy the homogeneous boundary conditions there exists a sequence  $\{w_i\}_{i=1}^{\infty}$  of functions such that  $(\{\phi_i\}, \{w_i\})$  is an  $L_0$ -optimal weighted residual pair.

Proof: Let  $w_i$  be the solution to  $Ay = \phi_i$  subject to the homogeneous boundary conditions. Existence is assured since  $\phi_i$  is admissable. We know that  $(y_0, \phi_j) = (y_n, \phi_j)$  so that

$$(Ay_n, w_j) = (y_n, Aw_j) = (y_0, Aw_j) = (Ay_0, w_j) = (f, w_j).$$

Thus  $(Ay_n, w_j) = (f, w_j)$  which is the system of equations for the method of weighted residuals. This proves  $(\{\phi_i\}, \{w_i\})$  is an  $L_2$ -optimal weighted residual pair.

Notice that if  $S_n(\phi)$  is invariant then any basis for  $S_n(\phi)$  and its image under A could be used as an optimal weighted residual pair. Also, in accordance with Theorem 4.5, any basis for  $S_n(\phi)$  is an  $L_2$  optimal Galerkin coordinate set.

Theorem 4.7. If  $A_1$  is strongly similar to  $A_2$  and  $\{\phi_i\}$  is an optimal Galerkin coordinate set with respect to  $A_1$  then  $\{P(\phi_i)\}$  is an optimal Galerkin coordinate set with respect to  $A_2$ .

Proof: By Theorem 4.5,  $S_n(\phi)$  is invariant under  $A_1$ . Since  $A_2$  is strongly similar to  $A_1$  we have

$$A_{2}P(S_{n}(\phi)) = PA_{1}(S_{n}(\phi)) = P(S_{n}(\phi))$$

or  $P(S_n(\phi))$  is invariant under  $A_2$  and again by Theorem 4.5 any basis for  $P(S_n(\phi))$  is an optimal Galerkin coordinate set with respect to  $A_2$ . This completes the proof.

Theorem 4.8. Suppose  $A_1$  and  $A_2$  are similar and that  $\{\{\phi_i\}, \{w_i\}\}$  is an optimal weighted residual pair with respect to  $A_1$ . Then  $\{\{\phi_i\}, \{P(w_i)\}\}$  is is an optimal weighted residual pair with respect to  $A_2$ .

Proof: We have  $A_1w_1=\phi_1$  and by similarity  $A_2P_2(w_1)=A_1(w_1)=\phi_1$ . Hence the result follows.

Theorem 4.9. Suppose  $A_1$  and  $A_2$  are similar. Then an almost optimal weighted residual pair with respect to  $A_1$  is an almost optimal weighted residual pair with respect to  $A_2$ .

Proof: Suppose ( $\{\phi_i\}$ ,  $\{w_i\}$ ) are almost optimal for  $A_1$ . Since  $A_1$  and  $A_2$  are similar

$$\| A_{2}(w_{i}) - \phi_{i} \| = \| P_{1}A_{1}(w_{i}) - \phi_{i} \| \leq \| P_{1}A_{1}(w_{i}) \| + \| \phi_{i} \|$$

Since  $P_1$  is continuous and linear it is bounded so that there is a K such that

$$\|P_1A_1(w_1)\| \le K\|A_1(w_1)\|.$$

Furthermore there is a number C such that

$$C > ||A_1(w_i) - \phi_i|| \ge |||A(w_i)|| - ||\phi_i|||$$

We can take the  $\{\phi_i\}$  such that  $||\phi_i|| = 1$ . Hence

$$||A(w_i)|| < C + 1$$

so that

$$\|P_1A_1(w_1)\| \le K(C+1)$$

and the first inequality above gives

$$||A_{2}(w_{i}) - \phi_{i}|| \le K (C + 1) + 1 \equiv M.$$

This completes the proof.

An additional and very significant result comes from using optimal coordinate and weight functions.

Theorem 4.10. If the optimal coordinate functions are used in the Galerkin method or an optimal pair is used in the method of weighted residuals then

the approximate solutions converge to the exact solution.

Proof: By optimality the approximate solutions are orthogonal projections of the exact solution on a sequence of monotonically increasing subspaces. Hence the approximate solutions converge to the exact solution.

From this theorem it is seen that if optimal coordinate and weight functions exist the method will converge.

Corollary to Theorem 4.6. If the conditions of Theorem 4.6 are satisfied the method of weighted residuals converges using the optimal pair.

Finally, it can be remarked that when optimal functions are used the sequence of norms  $\|y_n - y^*\|$  is nonincreasing. This follows since each  $y_n$  is the orthogonal projection  $S_n(\phi)$  and  $S_{n+1}(\phi) \supset S_n(\phi)$ .

### C. Improving Approximations

A description will now be given for the selection of an optimal weighted residual pair.

The usual procedure in the method of weighted residuals consists in selecting coordinate functions first and then weight functions. The choice of coordinate functions may be suggested by the problem but is usually a matter of taste. For example, if the solution is known to be periodic then of course periodic coordinate functions would be chosen. Frequently, however, the only constraint is that they should be linearly independent and satisfy the homogeneous boundary conditions. Once the

coordinate functions are selected the weight function choice is either based on the ease of resultant arithmetic (i.e., the computation of the inner products) or is completely arbitary. They should, of course, be linearly independent.

The following is suggested for these selections and is of course based on the preceding theorems. Select a sequence  $w_i$  of linearly independent functions that satisfy the homogeneous boundary conditions and such that there is enough "arbitrariness" (e.g., sufficiently many arbitrary constants) in each  $w_i$  so that  $Aw_i$  can be "forced" (e.g., the constants chosen) to satisfy the homogeneous boundary conditions.

It follows from Theorem 4.6 that the pair  $(\!\{\phi_i\}, \{w_i\}\!)$ , where  $\phi_i = Aw_i$ , form an optimal weighted residual pair. This method of choice is exhibited in the examples.

Another observation in connection with the above is quite helpful and can be used to obtain additional approximations.

Let

$$\overline{y}_n = \sum_{i=1}^n a_i (Aw_i)$$

denote the n<sup>th</sup> approximate solution. The system of equations to solve using  $\{Aw_i\}$  as coordinate functions and  $\{w_i\}$  as weight functions is

$$(Ay_n, w_j) = (f, w_j)$$
  $j = 1, ... n$ 

or

$$\sum_{i=1}^{n} a_{i} (A^{2}w_{i}, w_{j}) = (f, w_{j})$$

or

$$\sum_{i=1}^{n} a_{i} (Aw_{i}, Aw_{j}) = \sum_{i=1}^{n} a_{i} (\phi_{i}, \phi_{j}) = (f, w_{j}).$$

Interchanging the roles of weight and coordinate functions, now using

$$\bar{\bar{y}}_n = \sum_{i=1}^n b_i w_i$$

for the approximate solution, the system to solve is

$$(Ay_n, Aw_j) = (f, Aw_j)$$
  $j = 1, ... n$ 

or

$$\sum_{i=1}^{n} b_{i} (Aw_{i}, Aw_{j}) = \sum_{i=1}^{n} b_{i} (\phi_{i}, \phi_{j}) = (f, Aw_{j}).$$

Hence we conclude that the coefficient matrix to be inverted in each system is the same so that very little extra effort is needed to determine both  $\bar{y}_n$  and  $\bar{\bar{y}}_n$ , i.e., merely the additional calculation of the inner products on the right.

Once these two approximations are known they can be used to give an additional approximation of  $y^*$  as a linear combination of  $\bar{y}_n$  and  $\bar{\bar{y}}_n$  which involves only a 2  $\times$  2 matrix.

# D. Computational Considerations

In [45] a coordinate set is defined to be numerically optimal provided it is complete and the eigenvalues of the coefficient matrix for the n<sup>th</sup> approximation in the Galerkin method are bounded above and below independent of n. It is then shown that a numerically optimal set is invariant under similarity.

These ideas could easily be extended to the method of weighted residuals but have not been for two reasons. First, as shown in the examples, unless some prior information is available to give insight into the character of the problem, a random selection of coordinate functions may be just as good if not better than a set obtained using the idea of similarity.

Second, it is not reasonable to calculate eigenvalues of the coefficient matrix. The motivation for this idea is simply not to have the entries in the matrix get too large or too small and therefore create difficulties in machine calculations. This can usually be avoided by multiplication of each coordinate function by an appropriate constant. For example, if the coordinate functions were  $\sin n\pi x$  and the inner products in the coefficient matrix were of order n, and hence increasing greatly as the size of the matrix increased, by using  $\left\{\frac{\sin n\pi x}{\sqrt{n}}\right\}$  as coordinate functions, this difficulty would be overcome with no essential change in the problem. Therefore, the general rule is to properly "normalize" the coordinate and the weight functions to prevent the inner product entries in the coefficient matrix from becoming arbitrarily large or small.

It should be pointed out that the n<sup>th</sup> approximation may not be as good as the (n - 1)st approximation. That is, the sequence consisting of the norms of the difference between the n<sup>th</sup> approximation and the exact solution is not monotone decreasing. This can be observed in the examples. In practice, however, where the exact solution is not known there is in general no way to decide which particular approximation is better. All that is assured, if the method employed has been shown to

converge at all, in that for sufficiently large n the n<sup>th</sup> approximate solution can be made as close as desired in norm to the exact solution.

The following computational scheme is suggested by the preceding section and by observation of the examples.

Suppose several approximations have been obtained using  $\{\phi_i\}$  as coordinate functions. A sequential limit for the first few coefficients may be guessed giving, say, an  $n^{th}$  approximation  $y_n$  with some certainty. Then other approximations can be obtained using  $\{y_n, \phi_{n+1}, \phi_{n+2}, \dots\}$  as coordinate functions. This will necessitate inverting a much smaller matrix and hence a considerable saving in machine time.

## Chapter V

### A RELATED METHOD

The purpose of this chapter is to present a method for the solution of the type of problems under consideration, to prove that it is an optimal method and to prove that the method converges. The method is no more difficult to apply than those previously considered and has several reasons to be prefered.

Suppose throughout this chapter L is a linear operator defined in a separable Hilbert space and y\* is the unique solution to

$$L(u) = f (5.1)$$

subject to homogeneous boundary conditions. Let  $\{\phi_i\}$  and  $\{\psi_i\}$  be sets of linearly independent functions which satisfy the boundary conditions. Furthermore, let

$$u_{n} = \sum_{i=1}^{n} a_{i} \phi_{i}, \quad v_{n} = \sum_{i=1}^{n} b_{i} \psi_{i}.$$

<u>Definition 5.1.</u> The function  $u_n$  is called an <u>nth optimal approximation to the solution of 5.1 with respect to the set  $\{\phi_i\}$  provided</u>

$$\left\| L(u_n) - f \right\| \leq \left\| L \left( \sum_{i=1}^n c_i \phi_i \right) - f \right\|$$

for any choice of the coefficients ci.

Definition 5.2. The set  $\{\phi_i\}$  is a <u>better coordinate set for 5.1 than</u>  $\{\psi_i\}$  provided that if  $u_n$  and  $v_n$  are the respective  $n^{th}$  optimal approximations then

$$||L(u_n) - f|| \le ||L(v_n) - f||$$

These definitions are exactly the criteria commonly used, that is, we ask what function most nearly satisfies the operator equation.

The first result is a necessary and sufficient condition for a function to be an optimal approximation.

Lemma 5.1. The function  $u_n$  is the  $n^{th}$  optimal approximation with respect to  $\{\phi_i\}$  if and only if  $L(u_n)$  is the orthogonal projection of f on  $S_n(L(\phi))$ .

Proof: We have

$$- || L(u_n) - f || = || \sum_{i=1}^{n} a_i L(\phi_i) - f || \le || \sum_{i=1}^{n} c_i L(\phi_i) - f ||$$

for any choice of the c<sub>i</sub>'s. This is exactly the characterization of the orthogonal projection.

Lemma 5.2. If  $\{\overline{\phi}_i\}$  is an orthonormal basis for  $S_n(L(\phi))$  and  $u_n$  is the  $n^{th}$  optimal approximation with respect to  $\{\phi_i\}$  then

$$L(u_n) = \sum_{i=1}^{n} (f, \varphi_i) \overline{\varphi}_i. \qquad (5.2)$$

Proof: By Lemma 5.1,  $L(u_n)$  is the orthogonal projection of f on  $S_n(L(\phi))$  and truncation of the Fourier series gives the orthogonal projection.

Notice that equation 5.2 can be used to determine the approximation  $u_n$ . That is, having selected the  $\{\phi_i\}$  one can calculate  $\{L(\phi_i)\}$  and find (say by the Gram-Schmidt process) an orthonormal basis  $\{\overline{\phi}_i\}$  for  $S_n(L(\phi))$ .

Then equation 5.2 can be solved directly for the coefficients  $a_i$  in  $u_n$  using the linear independence of  $\{\overline{\phi}_i\}$  or a set of linear independent functions  $\{w_i\}$  can be used to give the system

$$\sum_{i=1}^{n} a_{i}(L(\varphi_{i}), w_{j}) = \sum_{i=1}^{n} (f, \overline{\varphi}_{i})(\overline{\varphi}_{i}, w_{j}) \quad j = 1, \dots n$$

for the determination of the a,'s.

Theorem 5.1. The set  $\{\phi_i\}$  is a better coordinate set for 5.1 than  $\{\psi_i\}$  provided the orthogonal projection of f on  $S_n(L(\phi))$  is a better approximation to f than the orthogonal projection of f on  $S_n(L(\psi))$ , that is

$$\left\| \sum_{i=1}^n \left( \mathbf{f}, \ \overline{\phi}_i \right) \ \overline{\phi}_i - \mathbf{f} \right\| < \left\| \sum_{1}^n \left( \mathbf{f}, \ \overline{\psi}_i \right) \ \overline{\psi}_i - \mathbf{f} \right\| \ .$$

Proof: By Lemma 5.2 the norm inequality becomes

$$||L(u_n) - f|| < ||L(v_n) - f||$$

which gives the result.

The question of convergence of this method is answered by the following result.

Corollary to Lemma 1: The norms  $||L(u_n) - f||$  converge monotonically to zero provided f is in  $\overline{S(L(\phi))}$ .

Proof: Since  $L(u_n)$  is the orthogonal projection of f on the subspaces  $S_n(L(\phi))$  and these increase in dimension it follows that the norms converge to zero. The monotonicty comes from the fact that, since

$$s_n(L(\phi)) \subset s_{n+1}(L(\phi))$$

the orthogonal projection of f on  $S_{n+1}(L(\phi))$  could be no worse an approximation to f than the orthogonal projection of f on  $S_n(L(\phi))$ .

To illustrate the application of the method presented in this chapter, we consider the equation

$$Ly = y'' + y = -x$$

which y(0) = y(1) = 0. As coordinate functions which satisfy the boundary conditions, we select

$$\varphi_n(x) = x^n(1-x)$$
 n = 1, 2, ....

To obtain an approximation in the form

$$y_2(x) = a_1 \phi_1(x) + a_2 \phi_2(x)$$

we determine an orthonormal basis  $\{\overline{\phi}_1,\ \overline{\phi}_2\}$  for  $S_2(L\phi)$ . An acceptable choice is

$$\overline{\varphi}_{1}(x) = \sqrt{\frac{30}{101}} \left(-2 + x - x^{2}\right)$$

$$\overline{\varphi}_{2}(x) = \sqrt{\frac{210}{2437}} \left(-6 + 13x - 3x^{2} + 2x^{3}\right)$$

We now require that

$$Ly_2 = (-x, \overline{\varphi}_1) \overline{\varphi}_1 + (-x, \overline{\varphi}_2) \overline{\varphi}_2$$

which is an equation involving the two unknowns  $a_1$  and  $a_2$ . By expressing the left-hand side as a linear combination of  $\overline{\phi}_1$  and  $\overline{\phi}_2$ , it

is possible to solve directly for a and a . However, if the inner product of each side with any two linearly independent functions is taken, two equations in a and a result, i.e.,

$$(\text{Ly}_2, \ \mathbf{w}_{\text{m}}) = \left(\sum_{p=1}^{2} \ (-\ \mathbf{x}, \ \overline{\phi}_{\text{p}}) \ \overline{\phi}_{\text{p}}, \ \mathbf{w}_{\text{m}}\right) \ \mathbf{m} = 1, \ 2.$$

We remark that if, in particular,  $w_m$  are chosen to be  $\overline{\phi}_m$ , the resulting system of equations are the same as would be obtained from the method of weighted residuals. For our example, to simplify computation we select

$$w_1 = 1, w_2 = x.$$

Then the system of equations, after simplification becomes

$$2a_1 + a_2 = \frac{55}{101}$$

$$a_2 = \frac{413}{2437}$$
.

It may be the case that the n<sup>th</sup> approximation obtained by this method is not the orthogonal projection of the solution on the span of the first n coordinate functions. However, we have the following result.

Theorem 5.2. If the operator L has a bounded linear inverse  $L^{-1}$  and if  $y_n$  is the  $n^{th}$  optimal approximation to the solution of

$$Ly = f$$

then

$$\lim_{n\to\infty} \|\mathbf{y}_n - \mathbf{y}\| = 0.$$

Proof: We have

$$\|y_n - y\| = \|L^{-1} (Ly_n - f)\| \|L^{-1}\| \|Ly_n - f\|$$

and since  $\| Ly_n - f \|$  can be made arbitrarily small the result follows.

### Chapter VI

## APPLICATIONS AND EXAMPLES

The purpose of this section is to apply several of the above methods to specific problems in order to illustrate the details of application, to demonstrate the approximation accuracy, and to point out certain difficulties which may arise.

We apply the Galerkin method and the method of weighted residuals to approximate the solutions to the following problems.

I. 
$$-y'' - y = x$$
 with boundary conditions  $y(0) = y(1) = 0$ 

II. 
$$y(x) - \int_{0}^{1} K(x, \xi) y(\xi) d\xi = \frac{1}{6} x (1 - x^{2})$$

where

$$K(x, \xi) = \begin{cases} x(1 - \xi) & \text{for } x \leq \xi \\ \xi(1 - x) & \text{for } \xi \leq x \end{cases}$$

III. 
$$((1 + x)y')' = x$$
 with boundary conditions  $y(0) = y(1) = 0$ 

In each of the problems considered, the operator is linear, selfadjoint, and positive definite. In addition, the operator in problem II is bounded, and hence continuous.

We consider two complete linearly independent sets of functions which satisfy the boundary conditions in each of the problems,

$$\phi_n(x) = x^n(1-x), \quad n = 1, 2, ...$$

$$\phi_{n}(x) = \sin n\pi x, \quad n=1, 2, ...$$

Thus, applying the general approximation theory of Chapter II, we take the space X to be the completion of the subspace of  $L_2$  spanned by  $\phi_n(x)$ ,  $n=1,\,2,\,\ldots$ 

To determine whether or not the Galerkin method will give the best  $L_2$  approximation using a fixed set of coordinate functions, we apply Theorem 4.5. That is, we determine whether or not the given operator applied to a linear combination of the first n coordinate functions gives an element which can also be written as such a linear combination. If the Galerkin method does not give the best approximation, applying Theorem 4.6, we determine a second linearly independ set of functions,  $w_n(x)$ ,  $n = 1, 2, \ldots$ , to be used as weight functions in the method of weighted residuals which will give the best approximation.

In each of the following examples, we find eight approximations, the third through the tenth. The  $n^{th}$  approximation to the solution of Ay = f is given by

$$y_{n}(x) = \sum_{j=1}^{n} a_{n,j} \phi_{j}(x)$$

We then evaluate the  $n^{th}$  approximation at each of the points  $x_i = 0.1(i-1)$ , i = 1, ..., 10 and compare these with the exact solution at each point. Finally, since the exact solution y is known, we are able to evaluate the norm of the difference of the exact and approximate solutions.

The three norms considered are defined as follows.

The L2-norm is

$$\|y - y_n\|_{L} = (y - y_n, y - y_n)^{\frac{1}{2}}$$

The A-norm is

$$||y - y_n|| = (Ay - y_n, y - y_n)^{\frac{1}{2}}$$
  
=  $(f - y_n, y - y_n)^{\frac{1}{2}}$ ,

i.e., the energy norm.

The sup-norm is

$$||y - y_n|| = \sup_{x \in [0,1]} |y(x) - y_n(x)|$$

In example 1, we use the Galerkin method to solve Problem I. We restate the problem equivalently as

$$y'' + y = -x$$
 with boundary conditions  $y(0) = y(1) = 0$ 

As coordinate functions which satisfy the boundary conditions, we select

$$\phi_n(x) = x^n(1-x), \quad n = 1, 2, \dots$$

Thus

$$A\phi_n = \phi_n'' + \phi_n = n(n-1)x^{n-2} - n(n+1)x^{n-1} + x^n - x^{n+1}$$

Applying the Galerkin method, we seek a  $j^{th}$  approximation in the form

$$y_j(x) = \sum_{n=1}^{J} a_n \phi_n(x)$$

where the j unknowns  $a_n$ , n=1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, \phi_{m}) = (-x, \phi_{m}), m = 1, ..., j$$

where

$$(A\phi_{n}, \phi_{m}) = \frac{m(m-1) + n(n-1)}{(m+n-1)(m+n)(m+n+1)} - \frac{m+n+\frac{1}{4}}{(m+n+2)(m+n+3)}$$

$$(-x, \phi_{m}) = -\frac{1}{(m+2)(m+3)}$$

The exact solution to Problem 1 is

$$y(x) = \frac{\sin x}{\sin 1} - x.$$

A(40)		-3,0441-53									•						Example 1.
A(9)	-2.3727-63	7,5758-03										•					X E
<b>9</b>	2.0752-03	-7.0125-03	Y10 [X]	00000	.019642	.036098	.051195	.062784	.069748	.071022	.165592	.052514	.030916	00000	6.391-06	5.350-05	1.399-35
(£) <b>v</b>	2.1973-45 -9.2773-03	1.1111-02	(X) 6A (I)	00000.	199611	5 .036094	4 .051185	5 .062766	4 .169722	5 .070987	5 .165554	779280. 2	988080. 8	000000. 0	5 1.941-05	1 7.284-05	3,197-05
A(6)	4.2725-04 -1.9531-03 2. 1.3889-02 -9.	-1.8219-02 1.	Y7 (x) YB (x)	000000. 000	542 .018641	104 .036655	301180. 825	160 . 162755	969690. 860	245 .070933	183 .165465	292280. 72	55 .030703	000000. 00	20-217:4-05	04 4.433-04	04 1.409-04
				000000	2 .018642	.036104	7 .051225	. 162860	268698. 5	3 .071245	. 165883	. 052827	. 031155	00000	1.840-04	9.624-04	3.250-04
¥(5)	7,7459-04 -6,1035-05 -9,7656-04 -1,0254-02	1,7857-02	<b>V6</b> (X)	000000	.019642	.03609	.051197	.062787	.069756	.071053	.065606	.052527	. 030923	000000	1,363-05	7,833-05	2,491-05
¥1.	-8.5912-83 -1.8124-82 -9.8942-83 -5.5711-83 -7.8198-83	-1.8250-02	<b>Y5</b> (X)	000000	.018642	.036098	.051195	.062784	.069749	.071021	.065588	.05250.	.030904	900000	1,742-06	7,745-06	2,875-06
A(3)	-2,3411-02 -1,0525-02 -9,4072-03 -9,8572-03 -1,0742-02	-7.6141-03	A (X)	000000.	.018641	.036090	.051195	.062783	.069746	.07101	.065585	.052503	. 930902	00000	3,696-07	5.630-06	5.895-07
A (2)	3-01 3-01 4-01 6-01 5-01	-0-	<b>13</b> (X)	00000.	.018625	.936105	.051220	.062803	.069746	.07099	.065561	.052496	.030918	0 0 0 0 0	1.658-15	2.060-06	2.565-05
		1.881	¥ (X)	000000-	.018642	.036098	.051195	.162783	. 069747	.07101	.165585	.152502	.330902	0	:28 <b>H</b>	34 34	SUP NORM
	1,8850-01 1,8850-01 1,8850-01 1,8850-01 1,8850-01	1.8840-01			-	?	r.	¥.	ĸ.	**	1-	•••	4	<u></u>	22 NORM	A NORM	8.75 6.75

In Example 2, we use the method of weighted residuals to solve Problem I.

As coordinate functions which satisfy the boundary conditions, again we select

$$\varphi_n(x) = x^n(1-x), n = 1, 2, \dots$$

As weight functions which will give the best approximations, we solve

$$Aw_n = \phi_n, n = 3, ..., 10$$

which gives

$$w_n = A_n \sin x - C_{n_1} \cos x + \sum_{k=1}^{n+2} C_{n_k} x^{k-1}$$

where

$$A_{n} = \frac{1}{\sin 1} \left( C_{n_{1}} \cos 1 - \sum_{k=1}^{n+2} C_{n_{k}} \right)$$

$$C_{n(n+2)} = -1, C_{n(N+1)} = 1 \text{ and}$$

$$C_{n_{k}} = -k(k+1) C_{n(k+2)}, k = 1, ..., n.$$

Applying the method of weighted residuals, we seek a j<sup>th</sup> approximation in the form

$$y_j(x) = \sum_{n=1}^{j} a_n \phi_n(x)$$

where the j unknowns  $a_n$ ,  $n=1,\ldots,j$  are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, w_{m}) = (-x, w_{m}) \quad m = 1, ..., j$$

$$\begin{aligned} &(\mathsf{A}\phi_{\mathtt{n}},\ \mathsf{w}_{\mathtt{m}}) \ = \ (\phi_{\mathtt{n}},\ \mathsf{A}\mathsf{w}_{\mathtt{m}}) \ = \ (\phi_{\mathtt{n}},\ \phi_{\mathtt{m}}) \ = \ \frac{1}{(\mathtt{m}+\mathtt{n}+\mathtt{1})(\mathtt{m}+\mathtt{n}+\mathtt{2})(\mathtt{m}+\mathtt{n}+\mathtt{3})} \\ &(-\mathtt{x},\ \mathsf{w}_{\mathtt{m}}) \ = \ -\sum_{k=1}^{\mathtt{m}+2} \ \frac{\mathsf{c}_{\mathtt{m}}}{\mathtt{k}+\mathtt{1}} \ + \ \mathsf{c}_{\mathtt{m}}_{\mathtt{1}} \left( \ \frac{\mathtt{1}}{\mathtt{sin}\ \mathtt{1}} \ - \ \mathtt{1} \right) \ -\sum_{k=1}^{\mathtt{m}+2} \mathsf{c}_{\mathtt{m}}_{\mathtt{k}} \left( \frac{\mathtt{cos}\ \mathtt{1}}{\mathtt{sin}\ \mathtt{1}} \ - \ \mathtt{1} \right) \end{aligned}$$

A:10:								-2.6154-02															Example 2.
A ( 9)							-2,1667-03	-5.7946-04															. Exa
F)						2. 4414-03	-4.8935-04	-4.8935-04	Y10 (X)	000000	159611.	.036085	.051184	.062709	.169752	.071004	.165569	.052516	.030907	000000	1,263-05	2.987-04	2.742-05
F. 37					1.9618-03	-4.8828-04 2	4.8828-03 -4	6.2408-02 -4	Y9 (X)	00000	.010642	. 136198	.051199	.062796	077690.	.071052	.165629	.052548	.030934	00000	2.624-05	1.270-04	4.650-05
				3-04					YD (X)	00000	.011642	.036101	.051199	.062786	. 169748	. 071 017	. 165583	.052505	.030915	00000	5.162-16	7.757-05	1.408-05
A (6)			20	03 0.3925-04	93 -0.7057-04	15 -1.787-64	3 -8.7837-04	12 -0.7957-04	Y7 (X)	00000	19841	. 036099	.051196	.062795	187690.	.071025	.065592	.052509	.030907	000000	4.242-86	2.264-05	7.080-06
A(5)			1.2207-03	-1.2372-03	-5.2032-03	-4. 3945-03	-1.2372-03	-5,5695-12	¥6 (x)	000000	.018642	.036098	961150.	.062786	.069752	.071024	165591	.052506	.030903	000000	3.396-06	1.300-05	5.841-06 7.080-06
71.63		-8.6960-03	-1.0986-02	-0.6060-03	-1.8232-03	-1.8232-03	-1.0376-02	-1.0232-03	Y5 (X)	. 00000	.018641	.036096	.051191	. 162777	.069738	.071005	.165569	.052407	.030891	000000	9.614-06	4.342-05	2,077-06 1,635-05
A (3)	-2.3436-12	-1, 0513-02	-9, 1553-13	-9.9030-03	-1.3603-02	-1, 4282-02	-7.8735-03	8.6212-03	<b>V&amp;</b> (X)	000000	.018641	.036098	. 051 195	.062785	.069747	.071020	.065507	. 152504	.030903	000000	1.009-06	7.227-06	2.077-06
(2) 4		1.8864-01 -1	1.8829-01 -9	1.8843-01 -9	1.8923-01 -1	1.8951-01 -1	1.8791-01 -7.	1.8128-01 8.	Y\$ [X]	000000	.018617	.036195	.051211	.062790	.069747	.071003	.165570	.152506	. 030925	00000	1.491-05	2.192-04	2.522-05
									¥ (X)	. 00000	.018642	. 036098	.051195	.062703	.069747	610170.	.065585	.052502	206050	00000	<b>.</b>	£	£ 2
161	1.8766-01	1,8857-01	1.8839-01	1. 8839-81	1,8834-01	1.8832-01	1. 6843-01	1.8903-01	* '		٦,	.2	۲,	٠.	ĸ.	۰.	*;	•	<b>.</b>		2 NORH	NO.	SUP NORM

In Example 3, we use the method of weighted residuals to solve Problem 1.

As coordinate functions which satisfy the boundary conditions, again we select

$$\varphi_n(x) = x^n(1-x), n = 1, 2, ....$$

As weight functions, we use

$$w_{n}(x) = \sin n\pi x, n = 1, 2, ...$$

Applying the method of weighted residuals, we seek a j<sup>th</sup> approximation in the form

$$y_j(x) = \sum_{n=1}^{j} a_n \phi_n(x)$$

where the j unknowns  $a_n$ , n = 1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, w_{m}) = (-x, w_{m}) \quad m = 1, ..., j$$

$$(A\phi_n, w_m) = (1 - (m\pi)^2) \int_0^1 x^n (1-x) \sin m\pi x dx$$

$$(-x, w_m) = \frac{1}{m\pi} (-1)^m.$$

																							•
A(10)								8,8319+03				•											Example 5.
A (9)							1.9175+13	1.9199+03															•
¥(0)						4. 5231 +02	4, 3198+62	-1.1427+04	710 (X)	000000	011132	6 150531	6 .554517	9 2.666458	5 4.401882	8643158	-,286620 -1,184796 -3,400172 -8,206575-16,846351	529025 -2.666043 -9.574152-30.274595	274248 -7.876163	5 .000015	10-561.1 0	1 1.718-02	11 3.033+01
A(7)					1.0107+02	1.0187+02	-2.2111+03	-2.2160+03 -	1 49 (X)	900000	96000 -	6 .192456	.534926	5 1,186598	510005	6 -2,86186	2 -8.20657	3 -9.57415		1 .000005	0 4.306+00	1 5.455+01	1.017+01
. •				=					<b>X9</b> (X)		.017407	112636	.262038	170615	593375	2.09076	3,60017	2.66604	279677.	100000.	1.631+00	1.016-01	3,545+00
¥.				2.5441+81	2.5441+01	-4.3676+02	-4.3644+12	3.5750+03	(X) LA	000000	.022975	.162230	.034677	177579	606775	- 600770.1	- 962981.1	529025 -	.622287	.00000	6.509-01	6.394+00	1.287+00
A(5)			6.8704+00	6.8705+00	-9.8022+01	-0.8021+01	5.4016+02	5.4098+02	4 <b>6</b> (X)	000000	.014809	957500"	065203	200344	349733	-,413963 -1,077009 -2,090766 -2,061868	286620 -	152536	380808	.00000	2.753-01	2.428+00	4.859-01 1.287+00
A(4)		2.0717+00	2.0717+00	-1.8100+01	-1.0101-01	7.6913+01	7.6863+81	-2.0230+02	<b>15</b> (x)	000000	.002144	021635	165102	107254	118933	074374	032207	.165567	.222241	00000	1.246-01	1.014+00	1,109-01 1,914-01
A (5)	7, 3638-01	7,3638-01	-5.7481+00	-3.7481+88 -	9.6303+00 -			-1.5204+01 -	YEIX	.00000	004846	014105	017615	006500	.023540	. 070148	.121741	.155023	.132502	. 000000	6.243-02	4.653-01	
									<b>73</b> (X)	. 60000	.006472	.016957	.037170	.063059	.190805	116821	127754	.120481	. 182115	000000	3,777-62	2.192-01	6.828-02
4(2)	3 5.4200-01	3 -7,1606-01	1 -7,1406-01	1 8.9257-01	2 0.9256-01	•			YIX	000000	.018642	036098	.051195	.162783	1969147	810160.	.165585	.052502	.030962	000000	£	-	<b>E</b>
	B. 1269-03	0.1278-03	1,2995-01	1,2995-01	5,6243-02	9, 6211-32	1,0038-01	1.0854-01	×	•	٠,	~:	s.	۳.	ĸ,	٠,	۲.	•	o,		22 NORM	+ NORM	S.P. NORM

In Example 4, we use the Galerkin method to solve Problem 1.

As coordinate functions which satisfy the boundary conditions, we select

$$\phi_{n}(x) = \sin n\pi x, n = 1, 2, ....$$

Thus

$$A\phi_n = \phi_n'' + \phi_n = \left(1 - (n\pi)^2\right) \sin n\pi x.$$

Applying the Galerkin method, we seek a j<sup>th</sup> approximation in the form

$$y_j(x) = \sum_{n=1}^{j} a_n \phi_n(x)$$

where the j unknowns  $a_n$ , n = 1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n} (A\phi_{n}, \phi_{m}) = (-x, \phi_{m}), m = 1, ..., j$$

$$(A\phi_{n}, \phi_{m}) = \begin{cases} \frac{1}{2} \left(1 - (m\pi)^{2}\right) & \text{if } m = n \\ 0 & \text{if } m \neq n \end{cases}$$

$$(-x, \phi_{m}) = \frac{1}{m\pi} (-1)^{m}.$$

Thus

$$a_n = \frac{2(-1)^n}{n\pi[1 - (n\pi)^2]}$$

and

$$y_j(x) = \frac{2}{\pi} \sum_{n=1}^{j} \frac{(-1)^n}{\pi [1 - (n\pi)^2]} \sin n\pi x.$$

A(19)								-6.4560-05															•	Example 4.
637							8.8592-05	8.8592-05																- E
¥(0)						-1.2610-04	-1.2619-04	-1.2618-04		EXI OSA	.00000	.018647	.036087	. 151211	.062760	.069778	778070.	.165641	.052423	.031019	00000	5.679-05	2.427-03	2.216-04
r					1.0045-04	1. 8845-04 -1.	1.0845-04 -1.	1.8865-66 -1.	,	Y9 (X)	.00000	.018647	.036007	.151211	.062760	.069770	776070.	.065641	.052423	.031019	00000	7.287-05	2.019-03	2.739-04
A(7)				ž						<b>19</b> (X)	.00000	.018619	. 136141	.051139	.162845	. 169689	.071961	.065570	.052475	.030991	. 000000	9.609-05	3.329-03	3,403-04
A(6)				-2.9947-94	-2.9947-04	-2.9947-84	-2.9947-04	-2.9947-54		Y7 (X)	000000	.018693	.036020	652150	.062770	.069689	.071136	.165450	.052595	716020.	000000	1.311-04	4.013-03	4.373-04
V(S)			5.1812-04	5.1812-84	5,1812-04	5.1812-06	5.1812-04	5.1012-04		76 (X)	000000	195610.	.136199	.051201	. 862668	.169877	.0711025	.065391	.052774	.030765	.00000	1.869-04	4.967-03	5.825-04
534		-1.0143-03	-1.0145-03	-1.0145-03	-1,0143-03	-1.0163-03	-1.0143-03	-1.0143-03		<b>45</b> (X)	000000	.018826	.036023	.051 025	. 062945	.069877	.070760	.065567	.052950	.030400	000000	2.825-04	6.368-03	90-151-04
A (3)	2.4162-03	2.4162-03	2.4162-03 -	2.4162-03	2.4162-03 -	2.4162-03 -	2.4162-03 -	2.4162-03 -		Y& (X)	000000	.010300	.036023	.051543	. 162945	. 169359	.070740	.066086	.05250	. 129962	.00000	4.626-04	8.576-03	1.214-03
<b>&amp;</b>	2724-03 2.	-8.2724-03 2.	-8.2724-15 2.	-8.2724-03 2.	-0.2724-03 2.	-8.2724-03 2.				<b>13</b> (X)	. 60000	.019272	.036619	.050947	161980	. 169359	.071705	.166682	.152354	.020997	. 0 0 0 0 0	8.535-04	1.242-02	2.000-03
7	Ŧ						12 -8.2724-03	02 -8.2724-03		YEX	000000	.018642	.036098	.051195	.862783	.069747	.071010	. 165585	.052502	.030902	.00000	£	_	### (F
4(1)	7,173-12	7,1775-02	7,1775-02	7,1775-02	7,1775-02	7,1775-02	7,1775-02	7,1775-02		×	7	٠,	7.	s.j	, <del>,</del> ,	ιν. ·	۰.	٠.	•	<b>.</b>	-	CZ NORM	A NORM	SUP NORM

In Example 5, we use the method of weighted residuals to solve Problem I.

As coordinate functions which satisfy the boundary conditions, we select

$$\phi_n(x) = \sin n\pi x, n = 1, 2, ...,$$

and as weight functions, we use

$$w_n(x) = x^n(1-x), n = 1, 2, ....$$

Applying the method of weighted residuals, we seek a  $j^{th}$  approximation in the form

$$y_j(x) = \sum_{n=1}^{j} a_n(A\phi_n, w_m) = (-x, w_m) m = 1, ..., j$$

$$(A\varphi_n, w_m) = \left(1 - (n\pi)^2\right) \int_0^1 x^m (1-x) \sin n\pi x dx$$

$$(-x, w_m) = \frac{1}{m\pi} (-1)^m.$$

•																						
A(10)							1, 3997+00				•											Example 5.
A(9)						-6.9335-01	-6.9916+00		,		•								•			•
A(B)					3.4994-01	3,1233+00	1.5206+01	K) 440 (K)	000000.	54 .003848	90 026912	69 .281528	.108250 -1.347255	16 4.852569	-5.107751 -6.687048	3.880348 -7.896488	92 45.171435	18-51.558666	500000- 20	10 2.160+01	02 4.770+02	10-537-01
A(7)				-1.7702-01	-1.4010+00	-6.0354+00	-1, 8934+01	YB (X) Y9 (X)		.00160 .00554	. 015174 . 050698	3680077069	.402301 .10820	0002 1.537306	1901 -5.1077	3,455134 3,88034	2.045007 11.416292 43.171435	4450-19.7644	20000000002	3+00 7,516+00	2+01 1.495+02	7.414+00 2.198+01
A ( 6)			8.6761-02	6.1781-01	2.3616+00	6, 9553+00	1.4940+01	Y (X) TY	00000	. 103521 . 100		.071244043680	.007101.	349731 500082	.175780 -1.421901	1.577915 3.45	031652 2.04	685092 -2.404034 -7.174450-19.764418-51.558666	00000 000001-	9.173-01 2.643+00	1.438+01 4.692+01	
A (S)		-3.7860-62	-2.5476-01	-8.7988-01	-2.1760+00	-4.4060+00	-7.8174+00	76 (x)	000000	000000	016476	.047230	024930	665010	. 383080	.530220 1.	- 607771	685092 -2.	000000	2.990-01 9.	4.104+00 1.	7,214-01 2,435+00
\$ F	9.3065-03	B. 5027-02	2.8717-01	6.5232-01	1.1902+00	1.0942+00	2,7525+00	<b>15</b> (X)	000000	.001955	917810.	.012033	.025412	.121078	.233317	.106701	011523	-,113600	00000	8.506-02	9.771-01	1.717-01
A (3)	6,1117-03	-5.9995-02	-1,4630-01	-2.5607-01	-3,7836-01	-5.0521-01	-6.3184-01	¥€ (X.)	000000	.001980	.015959	. 122225	.059574	.108050	.160248	.135055	. 196908	082730	000000	5.919-02	2.369-01	7,350-62
A (2)	-5.9832-02 (	-3.9279-02 -	-2.2819-02 -1	-7.1217-03 -2	6.6439-03 -3	1.8272-02 -5	2.7967-02 -6	<b>13</b> [X]	000000-	606000	. 008280	. 026702	.057304	.094897	127641	146510	.122088	.071326	00000	4.200-02	2.365-01	7.636-52
	1,0101-01 -5,5	9.8943-02 -3.9	9.7867-02 -2.2	9.7167-12 -7.1	9.6499-12 6.6	9.6131-02 1.8	5.5819-62 2.7	x Y(X)	000000	1 .018642	2 .036098	261138. 8	1 .062783	737690. 8	6 .071018	7 .165595	. 152512	5 .030902	63 63 63 63 63 63	22 NORM	A NORM	SUP NORM
A.01	5 6	<b>68</b> .	9.7	9. 19	29.6	5.61	5.58	•		٠,	7.	~;	٠,	ĸ.	7	•	-	. J.	-	3		ä

To conclude the analysis of Problem I, we compare the results of Examples 1 through 5.

A comparison of the first three examples using as coordinate functions

$$\phi_n(x) = x^n(1-x), n = 1, 2, ...$$

shows that the Galerkin method gives good results. In the second example where an optimal weighted residual pair is used, it can be seen that in most cases, the approximations are not much, if any, better than those obtained by the Galerkin method. In those cases where the Galerkin method apparently gives a better approximation, the discrepancy is attributed to the inaccuracy of numerical evaluation of the inner product expressions in the optimal weighted residual pair case and to the error produced in the numerical evaluation of the integrals defining the respective norms. This is substantiated by Theorem 4.10 since the sequence of norms  $\{||\mathbf{y}_{\mathbf{n}} - \mathbf{y}_{\mathbf{0}}||\}$  is non-increasing which is not reflected in this example. A comparison of examples 2 and 3 shows more clearly that the optimal weighted residual pair gives the best approximation using this set of coordinate functions.

In general, the problem of solving  $Aw_n = \phi_n$  to obtain an optimal weighted residual pair is as difficult as solving the original problem Ay = f. In Problem I, in particular, after obtaining the optimal weight functions, it was found that these functions were very close to being orthogonal to the right-hand side of the original equation. By

examination of the inner-product expression  $(f, w_n)$ , it can be seen that as n increases, we have  $(f, w_n)$  represented as the sum of three increasingly large factors to produce a term which tends to zero. Thus to obtain the accuracy required, it was necessary to evaluate  $(f, w_n)$  using double precision arithmetic.

In summary, we conclude that although the Galerkin method did not give the orthogonal projection of the solution on the span of the first n coordinate functions  $\phi_n(x) = x^n(1-x)$ , the increased difficulty of finding such a projection did not warrant the increased effort.

In Examples 4 and 5 coordinate functions

$$\phi_n(x) = \sin n\pi x, n = 1, 2, ...$$

were used. In this case

$$A\phi_n = [1 - (n\pi^2)] \sin n\pi x$$

which shows that the Galerkin method gives the orthogonal projection since, as in Theorem 4.5, the span of the first n coordinate functions is invariant under the operator A. Thus, as we would expect, the n<sup>th</sup> approximation using the Galerkin method in Example 4 is much better than the corresponding approximation obtained by the method of weighted residuals in Example 5.

As a final comparison we note that in Examples 2 and 4 orthogonal projections on the spans of the respective first n coordinate functions are obtained in each case. However, each approximation in Example 2 is

better than the corresponding approximation in Example 4. In fact, each approximation in Example 1, which is not the orthogonal projection, is also better than the corresponding approximation in Example 4. This indicates that the coordinate functions  $\phi_n(x) = x^n(1-x)$ ,  $n=1, 2, \ldots$  are a better set for Problem I than  $\phi_n(x) = \sin n\pi x$ ,  $n=1, 2, \ldots$  Although no criteria for the selection of a best set of coordinate functions is given, we remark that any additional information about the problem such as physical characteristics, periodicity of the solution, etc., are helpful in determining an optimal set of coordinate functions. In the case of Problem I, we speculate, for example, that coordinate functions  $\phi_n(x) = x^n(1-x)$  give a better approximation than  $\phi_n(x) = \sin n\pi x$  since the right-hand side of the equation, -x, is certainly not periodic.

In Example 6, we use the Galerkin method to solve Problem II.

As coordinate functions we select

$$\varphi_n(x) = x^n(1-x), n = 1, 2, ....$$

Thus

$$\begin{split} A\phi_{n} &= \phi_{n}(x) - \int_{0}^{1} K(x, \xi) \phi_{n}(\xi) d\xi \\ &= -\frac{2x}{(n+1)(n+2)(n+3)} + x^{n}(1-x) + \frac{x^{n+2}}{n+1)(n+2)} - \frac{x^{n+3}}{(n+2)(n+3)} . \end{split}$$

Applying the Galerkin method, we seek a jth approximation in the form

$$y_{j}(x) = \sum_{n=1}^{j} a_{n} \phi_{n}(x)$$

where the j unknowns  $a_n$ , n = 1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, \phi_{m}) = \frac{1}{6} (x - x^{3}, \phi_{m}), m = 1, ..., j$$

where

$$(A\phi_{n}, \phi_{m}) = \frac{2}{(m+n+1)(m+n+2)(m+n+3)} + \frac{2}{(n+1)(n+2)(n+3)} \times$$

$$\left[ \frac{2n+m+6}{(m+n+3)(m+n+4)(m+n+5)} - \frac{1}{(m+2)(m+3)} \right]$$

$$\frac{1}{6} (x - x^{3}, \phi_{m}) = \frac{2m+7}{3(m+2)(m+3)(m+4)(m+5)}$$

By means of Green's function, the boundary value problem, Problem I, can be transformed to the integral equation of problem II. Thus the solution to Problem II is also

$$y(x) = \frac{\sin x}{\sin 1} - x.$$

We observe here that, to within the accuracy of numerical methods, the n<sup>th</sup> approximation obtained by solving Problem II is the same as the corresponding approximation obtained by solving Problem I (Example 1). This fact is better illustrated by considering the coordinate functions

$$\phi_n(x) = \sin n\pi x, n = 1, 2, ....$$

Applying the Galerkin method, we seek a j<sup>th</sup> approximation in the form

$$y_{j}(x) = \sum_{n=1}^{j} a_{n} \varphi_{n}(x)$$

where the j unknowns  $a_n$ , n = 1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, \phi_{m}) = \frac{1}{6} (x - x^{3}, \phi_{m}) \quad m = 1, ..., j$$

where

$$(A\phi_n, \phi_m) = \begin{cases} \frac{1}{2} 1 - (m\pi)^2 & \text{if } m = n \\ 0 & \text{if } m \neq n \end{cases}$$

$$\frac{1}{6}(x - x^3, \varphi_m) = \frac{1}{m\pi} (-1)^m$$

Thus

$$a_n = \frac{2(-1)^n}{n\pi[1 - (n\pi)^2]}$$

which is the same as the result obtained in Example 4.

We conclude that all observations made concerning the operator and approximate method of determining a solution to Problem I are equally valid in determining a solution to Problem II.

•																							
A ( 10 )								-7,4760-04															Example 6.
¥19							-0.0501-04	-3.4688-04															· X
e e						-4, 1504-03	-4.4679-04	-4.4678-04	Y10 [X]		3 .010642	.03609	1051194	.062701	. 169743	. 1071011	.165572	.052483	.030882	.00000	1.021-05	7.746-05	2.002-05
A(7)					1.9684-03	6.0359-03 -	5.1128-03	4.6387-03 -	1 Y9 (X)	. 69660	1 .018643	660980. 7	161191	2 .062768	317690. 1	996010. 1	915591.	7 .052431	. 030850	.00000	4.001-05	2.071-04	7,263-05
A(6)				9.9182-04	-7.9966-04 1.	-7.9966-84 6.	-7.9966-04 5.	-7.9966-04	(X) A (X)	0000000	1 .018641	1636097	161191.	1 .162772	\$ .069724	\$700T0. 6	115591.	1 .052397	962080. 2	000000. 0	6 5,583-95	5 4.015-04	3.241-16 7.946-16 1.875-16 1.130-14
4			Ş						77 (x)	00000	. 118641	.036098	161194	. 162781	.069745	.07101	.165584	.052501	.030902	000000	1.021-06	1.192-05	1.875-0
A(5)			6.1635-04	-1.1245-03	-5.4092-03	-0.7891-03	-1,1245-83	-7.5989-63	Y6 (X)	000000	.018641	.036096	161191.	. 162777	. 169740	0.0170.	.165578	. 052499	.030902	000000	4,705-66	2,175-05	7.946-06
7		-8.6060-03	-1.0010-02	-0.6365-03	-1.6534-03	-1,6534-03	-9,1553-03	-1.6534-03	<b>15</b> (X)	900000	.018641	.036096	.051193	.062782	.069747	.071020	.065588	. 052505	.030903	.000000	1.807-06	1,105-05	
A (S)	-2.3637-02	-1.0529-02	-9. 5994-05	-1.0040-02	-1.5702-02	-1.2817-02	-1.0690-02	-1.2848-02	Y4 (X)	. 00000	.010641	.036098	.051195	.062783	.059747	.071018	.165586	.052503	.030902	. 00000	3.510-07	5.998-06	6.463-07
A (2)	1.9636-01 -2	1.0964-01 -1	1.0832-01 -9	1.8838-01 -1	1.0922-01 -1	1.8887-01 -1	1.8856-01 -1	1.8891-01	<b>YS</b> (X)	. 90000	.018617	.036095	.151211	. 162798	.069747	.071003	.165570	.052505	. 131925	000000	1.491-05	2.391-04	2.520-05
									YIXI	. 808080	.010642	.036098	.051195	.062783	.169747	.071018	.065585	.052502	.030902	000000	ž	<b>.</b>	- FE
ACT	1.0766-01	1.8837-81	1,0039-01	1.9840-01	1, 8834-01	1, 8657-01	1, 8840-01	1.6837-01	×	•	٠,	7	ri,	۳.	ĸ	٠.	.,	•	٠.	-	22 NORM	A NORM	SUP MORN

In Example 7, we use the Galerkin method to solve Problem III.

As coordinate functions which satisfy the boundary conditions, we select

$$\phi_n(x) = x^n(1-x), n = 1, 2, ....$$

Applying the Galerkin method, we seek a j<sup>th</sup> approximation in the form

$$y_j(x) = \sum_{n=1}^{j} a_n \phi_n(x)$$

where the j unknowns  $a_n$ , n=1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, \phi_{m}) = (x, \phi_{m}), m = 1, ..., j$$

where

$$(A\phi_{n}, \phi_{m}) = \int_{0}^{1} \left( (1 + x) \phi_{n}^{"} + \phi_{n}^{"} \right) \phi_{m} d_{x} = -\int_{0}^{1} (1 + x) \phi_{n}^{"} \phi_{m}^{"} d_{x}$$

$$= -\frac{mn}{m+n-1} + \frac{(mn + m + n)}{m+n} - \frac{(1 - mn)}{m+n+1} - \frac{(m + 1)(n + 1)}{m+n+2}$$

$$(x, \varphi_m) = \frac{1}{(m+2)(m+3)}$$

The exact solution to Problem III is

$$y(x) = \frac{1}{4} \left[ x(x - 2) + \frac{\ln (1 + x)}{\ln 2} \right].$$

4.898-42 1.966-42 -1.894-42 -1.894-42 -4.894-42 -1.894-43 -1.994-42 1.966-42 1.964-42 -2.364-42 -1.894-42 -2.364-42 -1.894-42 -2.364-42 -1.894-43 -1.184-41 -1.494-43 -1.364-42 -2.364-42 -1.364-42 -1.364-42 -1.364-42 -1.364-42 -1.364-42 -1.364-42 -1.364-42 -2.364-42 -1.364-42 -1.364-42 -1.364-42 -1.364-42 -1.364-42 -1.244-42 -2.364-42 -2.344-44 -2.344-44 -2.344-44 -2.344-44 -2.344-44 -2.344-44 -2.344-44 -2.344-44									
-5.999-12 1,963-12 -1,1064-12 -2,367-12 2,1264-13 -4,899-13 -4,899-12 3,4639-12 -2,367-12 2,1264-12 -1,1064-12 -2,1064-12	A(10)	-2.0752-02	•						Example 7.
-5.8959-12 1.9663-12 -1.8064-12 -6.2864-13 -6.959-12 -6.8959-12 2.9653-12 -2.9597-12 6.2864-13 -5.8959-15 -6.8957-12 -2.9597-12 -2.9	A(9)	4.5288-02				,			표
-6.6559-02 1.9663-02 -1.8864-02 -1.8864-03 -6.5864-03 -6.5559-02 2.1246-02 2.1872-01 1.8696-01 5.3596-01 -1.88641 -2.1877-02 5.9145-02 -5.9914-02 4.3654-02 1.3672-02 -4.3823-02 -4.3823-02 -1.2812-01 1.8696-01 5.3691-02 -7.1625-02 5.7107-02 5.9914-02 4.3654-02 -1.2012-01 1.8696-01 5.3691-02 -7.1625-02 5.7107-02 5.9914-02 4.3654-02 -1.2012-01 1.8696-01 5.3691-02 -7.1625-02 5.7107-02 5.9914-02 4.3654-02 -1.2012-01 1.8696-01 5.3691-02 -7.1625-02 5.7107-02 5.9914-02 4.3654-02 -1.2012-01 1.8618-01 5.3691-02 -1.2012-01 1.8618-01 5.3691-02 -1.2012-01 1.8618-01 5.3691-02 -1.2012-01 1.8618-01 5.3691-02 -1.2012-01 1.8618-01 5.3691-02 -1.2012-01 1.8618-01 5.3691-02 -1.82421 -	7(8)	. 6187-61 . 6187-61		•					2.023-04
-5.8989-12 1.9663-02 -6.8982-02 5.6458-02 -1.8864-02 -6.8982-02 5.1268-02 -1.8864-02 -6.8982-02 5.1268-02 -2.3697-32 6.2866-03 -6.8982-02 5.1268-02 -3.7720-02 2.1240-02 -5.8899-03 -6.1608-02 -2.2461-02 2.1873-01 -3.9258-01 3.3588-01 -7.1037-02 5.9143-02 -5.9814-02 4.1016-02 1.3672-02 -7.1037-02 5.9143-02 -5.9814-02 4.1016-02 1.3672-02 -7.1052-02 5.7107-02 -5.9814-02 4.1016-02 1.3672-02 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -5.9960-02 4.5654-02 -1.2207-03 -7.1052-02 5.7107-02 -0.0000 -0.000000 -0.00000 -0.00000 -0.00000 -0.00000 -0.00000 -0.00000 -0.000	ACT)A			•					5 9.473-05
-5.9999-02 1.9663-02 -6.8999-02 3.6438-02 -1.8096-02 -6.8992-02 3.6438-02 -1.8096-02 -6.8992-02 3.428-02 -2.9697-92 6.2066-03 -6.8992-02 5.1208-02 -3.7720-02 2.1240-02 -6.808-02 -2.2461-02 2.1875-01 -3.9258-01 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-02 -5.9914-02 4.1016-02 -7.1837-02 5.9143-030242419242429 -7.1837-02 5.9143-030242419242429 -7.1837-02 5.9143-030242419242429 -7.1837-02 5.9143-0302424192412429 -7.1837-02 5.9143-0302404800404890412609 -7.1837-0304126404048004048901600	5								4 9.218-05
-5.000-02 1.9663-02 -1.004-02 -6.6559-02 3.6450-02 -1.004-02 -2.5697-02 -6.0002-02 4.5555-02 -2.5697-02 -6.0002-02 5.1200-02 -2.1075-01 -6.0000-02 -2.2461-02 2.1075-01 -7.1057-02 5.9145-02 -5.9014-02 -7.1057-02 5.9145-02 -5.9014-02 -7.1057-02 5.9145-02 -5.9014-02 -7.1057-02 5.9145-02 -5.9014-02 -7.1057-02 5.9145-02 -7.1057-02 5.9145-02 -7.1057-02 5.9145-02 -7.1057-02 5.9145-02 -7.1057-02 5.9145-02 -7.1057-02 5.9145-02 -7.1057-03 -1.024241 -7.1057-03 -1.024254 -1.024257 -1.024254 -1.024257 -1.024254 -1.024254 -1.026019 -1.027999 -1.028019 -1		• • • • • • • • • • • • • • • • • • •		•					6 3.921-04
-5.090-02 1.9665-02 -6.6959-02 3.6450-02 -6.6959-02 3.6450-02 -2.2461-02 -2.2	A15			•					6 1.878-96
-5.8989-82 16.8882-82 56.9923-82 56.9923-82 56.9923-82 56.9923-82 57.1937-92 57.1937-92 57.1625-92 57.1625-92 57.1625-92 57.1625-92 57.1625-92 57.1625-92 57.1625-92 57.1625-92 57.1625-92 57.1625-92 61257 -7.1625-92 61257 -7.1625-92 61257 -7.1625-92 61257 -7.1625-92 61257 -7.16000000000000000000000000000000000000	A(4)	-1, 8864-8 -2, 3697-3 -3, 7721-0 2, 1875-0 -5, 9814-0 -5, 9860-0		•			6 0.225-0		90-009-1
-5.8989-82 -6.8882-82 -6.9923-82 -6.9923-82 -6.9923-82 -6.9923-82 -6.9923-82 -7.1937-92 -7.1937-92 -7.1937-92 -7.1937-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.1625-92 -7.16010 -010000	A (3)	1.965-62 5.6458-82 6.5353-82 5.1208-02 5.9145-02 5.9145-02 5.7007-02		•					5 5.015-06
	121		<b>43</b> (X)	•			2.138-0	3,151-00	3,452-65
		-1,4865-41 -5, -1,3954-81 -6, -1,3956-81 -6, -1,3959-81 -6, -1,3929-81 -6, -1,3929-81 -6,		•			LZ NORM	A NORT	SUP NORM

In Example 8, we use the Galerkin method to solve Problem III.

As coordinate functions which satisfy the boundary conditions, we select

$$\phi_n(x) = \sin n\pi x$$
,  $n = 1, 2, \dots$ 

Applying the Galerkin method, we seek a j<sup>th</sup> approximation in the form

$$y_{j}(x) = \sum_{n=1}^{j} a_{n} \phi_{n}(x)$$

where the j unknowns  $a_n$ ,  $n=1,\ldots,$  j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, \phi_{m}) = (x, \phi_{m}) m = 1, ..., j$$

where

$$(A\phi_{n}, \phi_{m}) = \begin{cases} - mn \left[ (-1)^{m+n} - 1 \right] \left[ \frac{m^{2} + n^{2}}{(m^{2} - n^{2})^{2}} \right] & \text{if } m \neq n \\ - \frac{3}{4} (m\pi)^{2} & \text{if } m = n \end{cases}$$

$$(x, \varphi_m) = -\frac{1}{m\pi} (-1)^m$$

					•															-		
A(18)							2,1853-05															Example 8.
A(9)						-5,1712-05	-4.8548-05							*								Exa
A ( 0 )					4.2106-05	3, 6042-05	3,4631-95	1 <b>416</b> (X)	000000.	3013148	1020227	032685	62980 9	041276	040462	036142	027967	016050	000000	2.650-05	1.470-03	8.957-05
Eja				-1, 1984-14	-1.0317-04	-1.0460-04	-1.0400-04	(X) 6A	000000. 01	7015165	12124224	600250 61	9 038626	9.21279	5040460	0036145	0 027966	1016052	000000 0	5 2.908-05	3 1.585-03	1.379-04 1.047-04
A(6)			9.7625-85	7.9611-05 -1	0.1030-05 -1	7,9525-05 -1,	7.9730-05 -1.	77 (x) <b>78</b> (x)		7115117 - 25117	218 024262	881 632839	659 138678	224041220	533040505	073036110	025027990	123016041	00000 000	-05 4.553-05	-03 2.027-03	
A(5)		-5, 11 96-94		-2.8890-04 7	-2.8757-04 0	-2.8015-04 7.	-2.8788-04 7.	76 (X)	000000	813035 813152	024324 024218	03288032881	138578138659	041550041224	040403040533	036030036073	028122028025	015951016023	000000000000	9.065-05 5.119-05	3.057-03 2.227-03	2.364-04 1.638-04
A(6)	1420-04			2.6186-04 -2	2.6230-04 -2	2.6147-04 -2	2.6163-04 -2	<b>75</b> (X.)	. 000000	115154	824278	067220	038672	041340	040395	0360770	028182	0158750	000000	1.062-04 9.0	3,465-03 3.0	2.989-04 2.3
A (S)	-1.4154-03			-1.3649-03	-1,3640-03	-1.3645-03 2	-1.3641-13 2	<b>46</b> (X)	. 00000	012700	024215	033096	0.128210	041049	040352	-, 036353	020213	015621	000000-	2.530-04 1.062-04	5.401-03	4.957-04
1214	1.8922-03 -			1.0091-03 -1	1.8896-13 -1	1.0865-03 -1	1.8887-03 -1	<b>75</b> (X)	000000	013146	024489	032968	038413	041019	040637	036567	028088	015370	00000	2.972-04	6.539-13	6.949-04
	-6.2636-02 1.			-4,2410-02 1.	-4.2438-02 1.1	-6.2609-02 1.1	-4.2408-02 1.1	X Y(X)	000000. 0.	-1013124	.2 024261	.5032072	.4038643	.5141259	.6 848482	.7036116	.9928001	.9016000	000000 0.5	C2 NORTH	A NORM	SUP NORM

In Example 9 we use the method of weighted residuals to solve Problem III.

As coordinate functions which satisfy the boundary conditions, we select

$$\varphi_n(x) = x^n(1-x), n = 1, 2, ....$$

It can be shown that the operator in Problem III,  $A_1y = ((1 + x)y')'$  is similar (Definition 4.5) to the operator  $A_2y = y''$ . For weight functions, we solve

$$A_2 w_n = \varphi_n$$
  $n = 1, 2, ...$ 

to obtain

$$w_n(x) = \frac{x}{(n+1)(n+2)(n+3)} \left\{ x^{n+1} \left[ (n+3) - (n+1)x \right] - 2 \right\}.$$

Thus, applying Theorem 4.9, since  $\left\{\phi_n,\,w_n\right\}_{n=1}^{\infty}$  is an optimal weighted residual pair with respect to the operator  $A_2$ ,  $\left\{\phi_n,\,w_n\right\}_{n=1}^{\infty}$  is an almost optimal weighted residual pair with respect to the operator  $A_1$ .

Applying the method of weighted residuals, we seek a j<sup>th</sup> approximation in the form

$$y_{j}(x) = \sum_{n=1}^{j} a_{n} \phi_{n}(x)$$

where the j unknowns,  $a_n$ , n = 1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, w_{m}) = (-x, w_{m}) \quad m = 1, ..., j$$

where

$$(A\phi_{n}, w_{m}) = \frac{1}{(m+1)(m+2)(m+3)} \left\{ (m+3) \left[ \frac{n(m+2)}{m+n+1} - \frac{(mn+n+m+2)}{m+n+2} - \frac{(mn+2n+1)}{m+n+2} + \frac{(m+1)(n+1)}{m+n+4} \right] + \frac{2}{(n+1)(n+2)} \right\}$$

$$(x, w_{m}) = -\frac{(2m^{2} + 9m + 7)}{3(m+1)(m+2)(m+3)(m+4)(m+5)} .$$

4(10)				2.0930+00														Example 9.
¥(9)			3. 3626+00	3,4688-04														Exa
A(B)			1,1167-01	4.4678-04	Y18 (X)			039474	042155	141 757	036850	029541	021185	008112	000000	8.278-03	9.363-12	1.675-02
A(7)		-4, 6421-01	-1.7448-01 1	-4.0054+00	₹		2596625	039508	042229	041694	036756	029594	021186	000003	. 60000	8.204-03	9.320-02	1.677-02
A(6)		_	7.9964-64 -1.	7.9964-01 -4.	<b>YB</b> (X)		3 129 /9 /	1 139452	7042086	1041700	036818	7029568	5021100	000000	000000	8 0.272-03	9.170-02	1.673-02
¥					Y7 (X)		020054	04040	040807	040534	037779	030827	019796	007905	00000	0.130-03	7.667-02	1.586-02
	-1.5682+88		9.5014+00	1,1245-03	Y6 (X)		037690	039992	041007	040090	037865	030564	019627	000019	00000	0.105-03	7.849-02	1.617-02
ACC	1,4484+00 4,6326+00	2.8334+00	7.9890+00	7.2525+00	<b>YS</b> (X)		037520	039864	041112	041072	037905	030481	019614	008162	00000-	8,191-03	7.914-02	1.629-02
A (3)	-1,8749-01 -2,4021+00 -4,6138+00	-3.6057+00	-5.0255+00	-7.4064+00	<b>46</b> (X)		038059	-,040505	640949	040316	037559	031002	020155	007455	00000	8.144-03	7.809-02	1.575-02
15) 4	5,6527-01 -1 1,5719-00 -2 1,9754-00 -4		2.3933+00 -5	2.7537+00 -7	<b>13</b> (X)		123676	045137	045751	041690	035294	020572	020394	011020	00000	7.636-03	5.808-02	1,297-92
					. 4(x)		02424	032072	038643	141259	040482	136116	028001	016000	000000-	L2 NORM	F.	SUP NORM
A193	-2,9752-01	-4,6389-01	-5.2112-01	-5, 6064-91	×	<b></b>	. 2	s,	<b>*</b>	s;	•	۲.	•	٠.		3	A NORM	d.:S

A comparison of the results in Examples 7 and 9 shows that in all of the norms considered, the Galerkin method gives better approximations than the method of weighted residuals using an almost optimal weighted residual pair. Thus we are lead to conclude that in selection of the best approximation method, unless a priori information is available, the Galerkin method will give relatively good results.

To illustrate the use of the Galerkin method for partial differential equations, we consider the problem

IV. 
$$\nabla^2 u + x \frac{\partial u}{\partial x} = 2x^2 + 2y^2 + 2x^2y^2 - 2x^2y - xy^2 + xy - 2x - 2y$$
  
with boundary conditions  
 $u(x, 0) = u(x, 1) = u(0, y) = u(1, y) = 0.$ 

As coordinate functions which satisfy the boundary conditions, we select

$$\varphi_n(x) = \sin p\pi x \sin q\pi y$$

where the correspondence of p and q to n is given by

n	1	2	3	4	5	6	7	8	9	10
q	1	2	1	3	2	1	4	3	2	1
q	1	1	2	1	2	3	1	2	3	4

$$A\phi_n = \nabla^2 \phi_n + x \frac{\partial u}{\partial x} = \left[ -\pi^2 (p^2 3 q^2) \sin p\pi x + p + x \cos p\pi x \right] \sin q\pi y.$$

Applying the Galerkin method, we seek a  $j^{th}$  approximation in the form

$$u_{j}(x) = \sum_{n=1}^{j} a_{n} \phi_{n}(x)$$

where the j unknowns  $a_n$ , n = 1, ..., j are determined from the following system of j equations

$$\sum_{n=1}^{j} a_{n}(A\phi_{n}, \phi_{m}) = (f, \phi_{m}) \quad m = 1, ..., j$$

where f is the right-hand side of equation IV and if  $\phi_m(x) = \sin \, r\pi x \, \sin \, s\pi y \text{,}$ 

$$(A\phi_{n}, \phi_{m}) = \begin{cases} \frac{pr}{2(p^{2}-r^{2})} (-1)^{p+r} & \text{if } q = s, p \neq r \\ -\frac{\pi^{2}}{4} (p^{2}+q^{2}) - \frac{1}{8} & \text{if } q = s, p = r \\ 0 & \text{if } q \neq s \end{cases}$$

$$(f, \varphi_{m}) = \frac{1}{(s\pi)^{3}} \left[ 2A(r) - (s\pi)^{2} B(r) \right] \left[ (-1)^{s} - 1 \right]$$

where

$$A(r) = -\frac{1}{r\pi} (-1)^{r} + \left[ \frac{4}{(r\pi)^{3}} - \frac{2}{(r\pi)} \right] \left[ (-1)^{r} - 1 \right]$$

$$B(r) = \frac{4}{(r\pi)^{3}} \left[ (-1)^{r} - 1 \right].$$

The solution to Problem IV is

$$u(x, y) = xy(1 - x)(1 - y).$$

We remark that this problem was considered by Mihklin [38]. He applied the Galerkin method as we have done, however, his results are incorrect due to sign errors in his calculations of  $(A\phi_n, \phi_m)$  and  $(f, \phi_m)$ .

0.7	.214	¥(3)	ĀTĒ	115	¥(6)	EJY	A(B)	¥(9)	A140)	
6, 6238-12	6,6238-02 -5,4195-05 0,0000	0000 0							•	
6,6118-52	6.6118-32 -3.2460-03	0000.6	3,5338-03							
6.6118-02	6.6118-02 -3.2460-03	0.00.0	3,5338-03	0.0000						
6,6118-02	6.6118-02 -3.2660-03 0.0000	0.000.0	3,5330-03	0000	2.4640-03					
6.6105-02	6.6135-02 -3.2318-03	0.000.0	5.5144-03	0.000.0	2.4649-05	2,4640-03 -5,5073-04				
6.6135-02	6.6135-02 -3.2318-03	0.0000	3.5144-03	0.0000	2.4640-03	2.4640-03 -5.5073-04	0.0000			
6,6105-02	6.6105-02 -3.2318-03	0.000.0	3.5144-03	6.0000	2.4641-03	2.4641-03 -5.5073-04	0.000.	-4, 91 25-05		
6.6135-12	6.6135-02 -3.2318-03 0.0000	0.000.0	3.5144-03	0.00.0	2.4641-03	2.4641-03 -5.5073-04 0.0000	0.000.0	-4.9123-05	0.000	

Example 10.

As a final example we consider the application of the Galerkin method to the initial value problem

V. 
$$y'' + xy = 1$$
 with initial conditions  $y(0) = y'(0) = 0$ .

In attempting to apply the Galerkin method, we are immediately faced with two problems.

The first problem arises in the selection of coordinate functions which satisfy the initial condition. It is not difficult to satisfy the initial conditions, for

$$\phi_n(x) = x^{n+1}, \quad n = 1, 2, ...$$

is one such set. However, the set

$$\varphi_n(x) = x^{n+1}(R - x), n = 1, 2, ...$$

for any real number R, also satisfies the initial conditions. Thus, by an appropriate choice of R, we are able to force our approximations to vanish at any desired point.

Assuming an appropriate set of coordinate functions has been selected, the second problem arises in the calculation of inner products. Since

$$(f, g) = \int_a^b fgdx$$

we must have some finite interval [a, b] over which the integration is to be performed. We solve this problem by assuming that an approximate solution is desired only for certain values of the independent variable. Thus we may assign to b, any value such that [a, b] contains all value in which we are interested. In [31], Ames suggests that b be chosen as some point at which steady-state is achieved (see discussion in Chapter III). We point out, however, that this is not always possible. In Problem V, in particular, noting that the solution is

$$y(x) = x^2 \sum_{n=0}^{\infty} (-1)^n \frac{n!}{(3n+2)!} (3x^3)^n$$

steady-state is never achieved; in fact, the solution becomes increasingly unstable as x increases.

We conclude that for Problem V the Galerkin method and the method of weighted residuals are not applicable unless a second zero of the function y is known, in which case, we approximate the solution as shown in the first three boundary value problems.

## Chapter VII

### THE NONLINEAR CASE

The methods under consideration can, in principle, be applied to nonlinear problems as noted in Chapter III. However, there are at least two difficulties involved. First, and less serious although by no means minor, is the fact that the system of equations obtained for the coefficients is nonlinear. Hence it may be an impossible task to solve them. Second, there are no known convergence proofs for this class of problems.

The above drawbacks have not prevented the use of a method. The first has been overcome by the use of very few terms and the second has been ignored. Authors invariably claim "good" results.

The following is a summary of the use of Galerkin's method in a nonlinear situation.

Cesari [33], in a very elegant paper, uses a Galerkin approach to give an existence proof for a solution to certain nonlinear equations.

As an example he analyzes the problem

$$x'' + x + \alpha x^3 = \beta t$$
,  $0 \le t \le 1$ 

with boundary conditions

$$x(0) = 0, x'(1) + hx(1) = 0$$

where  $\alpha$ ,  $\beta$  and h are numbers.

Following in this direction Urabe [42] considers the nonlinear periodic equation

$$\frac{dx}{dt} = X(x, t)$$

where x and X(x, t) are vectors of the same dimension and X(x, t) is periodic of period  $2\pi$  in t. In particular he proves the following theorem.

Theorem: The existence of an isolated periodic solution,  $\hat{x}$ , of (7.1) implies the existence of Galerkin approximations  $x_m$  for sufficiently large m. The existence of Galerkin approximations  $x_m$  for sufficiently large m implies the existence of an exact solution provided a certain boundedness condition holds.

In a later publication Urabe [44] compares his earlier technique with a method of averaging when  $X(x, t) = \lambda \overline{Y}(x, t)$  where  $\lambda$  is a small parameter. A somewhat more detailed numerical analysis is given in [43].

Locker [50] generalizes the work of Cesari to give an existence analysis of

$$Lx = Nx$$

where L is an unbounded linear operator and N is a nonlinear operator. When L is self-adjoint his results reduce to Cesari's. He also gives a simple example.

There is one other recent reference which mentions Galerkin's method in connection with nonlinear problems. Bellman [32] considers a method for replacing a nonlinear problem by a linear one. As an alternative to the main point of his paper, he briefly indicates that a one term Galerkin approximation could be used to implement his method.

In summary, there have been two uses of Galerkin's method on nonlinear problems. First, as indicated above, a few existence results have been obtained for certain problems by employing a Galerkin approach. Second, and at the other extreme, the method has been applied without regard to convergence.

To illustrate some of the difficulties which arise is attempting to apply the Galerkin method or method of weighted residuals to non-linear equations, we consider the following differential equation.

$$(y')^2 + \pi^2 y^2 = \pi^2$$
 with boundary conditions  $y(0) = y(1) = 0$ 

We note first that  $y(x) = \sin \pi x$  is a solution to the equation, but that it is not unique, for  $y(x) = -\sin \pi x$  is also a solution.

In applying the method of weighted residuals, as coordinate functions which satisfy the boundary conditions, we select

$$\phi_n(x) = x^n(1-x), n = 1, 2, ...$$

and as weight functions, we use

$$w_n(x) = x^{n-1}, n = 1, 2, ....$$

We seek a j<sup>th</sup> approximation in the form

$$y_{j}(x) = \sum_{n=1}^{j} a_{n} \varphi_{n}(x)$$

where the j unknowns  $a_n$ , n = 1, ..., j are determined from the following system of j equations

$$(A \sum_{n=1}^{j} a_n \phi_n), w_m) = (\pi^2, w_m), m = 1, ..., j.$$
 (7.2)

Thus we obtain a system of non-linear equations which may be difficult or impossible to solve.

To illustrate, consider the case where j = 2. The system of equations (7.2) becomes

$$\left(\frac{\pi^2 + 10}{30}\right) a_1^2 - \left(\frac{\pi^2 - 10}{30}\right) a_1 a_2 + \left(\frac{71\pi^2 + 14}{105}\right) a_2^2 = \pi^2$$

$$\left(\frac{\pi^2 + 10}{60}\right) a_1^2 - \left(\frac{29\pi^2 - 49}{210}\right) a_1 a_2 + \left(\frac{485\pi^2 + 84}{840}\right) a_2^2 = \frac{\pi^2}{2}$$

Solving this system, assuming  $a_2 \neq 0$ , we obtain

$$a_1^2 = \frac{\pi^2}{\left[y(\pi^2 + 10) - 7(\pi^2 - 10) k + 2 (71\pi^2 + 14) k^2\right]}$$

$$a_2 = ka_1, \text{ where } k = \frac{2(51\pi^2 - 28)}{201\pi^2 + 28}.$$

Thus, in seeking a two-term approximation, we find two possible solutions. Continuing in this manner, it can be seen that as the number of terms in the approximation increases, the number of possible solutions to the system of equations (7.2) also increases. It is apparent that without a priori information regarding the solution of a non-linear problem, a purely arbitrary choice among the possible solutions may not give a desirable approximation.

### Chapter VIII

# APPLICATIONS AND ANALYSES OF METHODS OF THE GALERKIN TYPE AVAILABLE IN THE LITERATURE

This chapter consists of a rather complete representative cross section of references to the methods under consideration. Following each reference is a remark to indicate how the article uses a method. Some insight into typical applications, frequence of use and care in application can be gained from these references.

The list does not include many of the very brief references to one or more of the various methods. A reference is included only in case there is a fairly complete discussion of a method or a nontrivial application. Inclusion of a reference is not an endorsement of the application. Convergence proofs or justification for use almost never occur.

The references are grouped by (A) Ordinary Differential Equations, (B) Partial Differential Equations and (C) Studies of Methods per se.

### A. Ordinary Differential Equations

1. Birikh, "On the spectrum of small perturbations of plane-parallel Couette flow", (PMM) Journal of Applied Mathematics and Mechanics, Vol. 29, 1965, 946-949.

Solves a fourth order ordinary differential equation subject to homogeneous boundary conditions using the coordinate functions

$$V_{n} = K_{n} \left[ \frac{\cosh ax}{\cosh a} - \frac{\cos \sqrt{\lambda_{n} - a^{2} x}}{\cos \sqrt{\lambda_{n} - a^{2}}} \right]$$

for even n and the same funtions with sinh and sin replacing cosh and cos, respectively, for odd n. The same problem is solved in [9] using different coordinate functions.

2. Bruslinskaja, "Limit cycles for equations of motion of a rigid body and Galerkin equations for hydrodynamics", <u>Soviet Mathematics</u>, Vol. 5, Aug. 1964, 1051-1054.

Solves a first order system of three equations with little comment.

3. Chandraskhar, "Hydrodynamic and hydromagnetic stability", Clarendon Press, Oxford, England.

Solves the system

$$(D^2 - a^2)^2 u = (1 + ax) v$$
  
 $(D^2 - a^2) v = -\lambda u$ 

subject to u = Du = v = 0 for x = 0, 1 using as coordinate functions  $v_n$  =  $\sin n\pi x$  and  $u_n$  the solutions of

$$(D^2 - a^2) u_n = (1 + ax) \sin n\pi x$$

$$u_n = Du_n = 0 \text{ for } x = 0, 1$$

and using  $\sin n\pi x$  as weight functions.

4. DiPrima, "Applications of Galerkins method to a problem in hydro-dynamic stability," Quarterly of Applied Mathematics, Vol. 13, 1, 1955, 55-62.

Solves a problem similar to that in [3] using e<sup>nx</sup> times a polynomial as coordinate functions.

5. DiPrima, Walowit, Tsao, "Stability of flow between arbitrarily spaced concentric cylindrical surfaces including the effect of a radial temperature gradient," <u>Journal of Applied Mechanics</u>, Vol. 31, 1964, 585-593.

Solves a 2 × 2 homogeneous system of ordinary differential equations with homogeneous boundary conditions using  $(x^2 - \frac{1}{4})^2 \ x^{n-1}$  as coordinate functions for both unknowns. Shows no work.

6. DiPrima, "Some variational principles for problems in hydrodynamic and hydromagnetic stability", Quarterly of Applied Mathematics, Vol. 18, 1961, 375-385.

Mentions the application of the Galerkin method to a problem similar to [3]. He uses cosines as one set of coordinate functions and those of [10] for the other.

7. Duncan, "Torsional oscillation of a cantilever when the stiffness is of composite origin," TR Aero Research Commission, Great Britian Air Ministry, Part 1, 1937, 1809, 471-483.

Solves a second order ordinary differential equation with homogeneous boundary conditions using polynomials as coordinate functions.

8. Frazer, Jones and Shan, "Approximation to functions and the solution of differential equations", TR Aereo Research Commission, Great Britian Air Ministry, Vol. 1, 1937, 1799, 517-549.

Solves ordinary differential equations of order less than four using certain polynomials as coordinate functions.

9. Gallagher and Mercer, "On the behavior of small disturbances in plane couette flow", <u>Journal of Fluid Mechanics</u>, Vol. 13, 1962, 91-100.

Solves a fourth order equation with homogeneous boundary conditions using coordinate functions similar to those in [10].

10. Harris and Reid, "On orthogonal functions which satisfy four boundary conditions - Tables", <u>Astrophysical Journal Supplement</u>, Vol. 3, 1958, 429-452.

An orthogonal sequence of functions is presented which satisfy

$$y^{(3)} = ay, y = y^{1} = 0 \text{ at } x = \pm 0.5$$

$$c_{m}(x) = \frac{\cosh \lambda_{m} x}{\cosh \frac{1}{2} \lambda_{m}} - \frac{\cos \lambda_{m} x}{\cos \frac{1}{2} \lambda_{m}}$$

$$S_{m}(x) = \frac{\sinh \mu_{m} x}{\sinh \frac{1}{2} \mu_{m}} - \frac{\sin \mu_{m} x}{\sin \frac{1}{2} \mu_{m}}$$

where  $\lambda_{\,\,\mathrm{m}}^{}$  and  $\mu_{\,\,\mathrm{m}}^{}$  are the positive roots of

$$\tanh \frac{1}{2}\lambda + \tan \frac{1}{2}\lambda = 0$$

$$\coth \frac{1}{2}\mu - \cot \frac{1}{2}\mu = 0.$$

11. Krueger and DiPrima, "The stability of a viscous fluid between rotating cylinders with axial flow", <u>Journal of Fluid Mechanics</u>, Vol. 19, 1964, 528-538.

Solves a complex eigenvalue problem for a 2  $\times$  2 system using Galerkins method with the coordinate functions of [10].

12. Kurzweg, "Magnetohydrodynamic stability of curved viscous flows", TRII-29, Princeton University, Department of Physics, 1961.

Applies the Galerkin method to three simultaneous fourth order linear ordinary differential equations with homogeneous boundary conditions using polynomial coordinate functions.

13. Penzes and Burgin, "Free vibrations of thin isotropic oblatespheroidal shells", General Dynamics/Convair, San Diego, Calif.

Solves a second order homogeneous ordinary linear differential equation with variable coefficients using associated

Legerdre functions with the Galerkin method where the constraints are that the solution be single valued and bounded.

14. Rao, "The fundamental flexural vibrations of a cantilever beam of rectangular cross section with uniform taper", Aeronautical Quarterly, May 1965, Vol. 16, 139-144.

Applies Galerkin method to a fourth order ordinary equation with homogeneous boundary conditions.

15. Sastry, "Influence of wall conductance on the stability of Hartmann flow in a curved channel", Physical Society of Japan Journal,

Vol. 21, 6, June 1966.

Applies Galerkin method to an eigenvalue problem involving a fourth order differential equation.

16. Sennet and Wait, "Analog computer analysis of static beams", AICE,

Engineering Mechanics Division Journal, Vol. 91, April (I)

1965, 47-62.

Compares a Galerkin solution to an analog computer solution of a fourth order equation and finds the Galerkin solution within 3% of the other.

17. Smith, "On the growth of Taylor-Garther vortices along highly concave walls", Quarterly of Applied Mathematics, Vol. 13, 1955, 233-262.

Applies Galerkin method to two simultaneous ordinary differential equations, one of third order and one of fourth, with homogeneous boundary conditions using e<sup>-an</sup> times polynomials as coordinate functions.

### B. Partial Differential Equations

18. Bickley, "Experiments in approximating solutions of a partial differential equation", Philosophical Magazine, (7) 32, 1941, 50-56.

> Solves the one dimensional wave equation with certain boundary and initial conditions using his modification of the Galerkin method, collocation and least squares obtaining one and two term approximations.

19. Biezeno, Grammil, Blackie, "Theory of Elasticity", Engineering

Dynamics, Vol. 1, 170-176.

Applies Galerkin method to  $\nabla^4 u = f$  with homogeneous side conditions.

20. Bolotin, "Nonconservative Problems of the Theory of Elastic Stability", 1963, pp. 58-62, 108-109, 207-212, 247-265.

Makes several applications of the Galerkin method.

21. Cheng and Pan, "Stability analysis of gas lubricated, self-acting, plain, cylindrical journal bearings of finite length using Galerkins method", <u>Journal of Basic Engineering</u>, Vol. 87, 1965, 185-191.

Applies Galerkin method to a second order equation with homogeneous boundary conditions using cosines as coordinate functions.

- 22. Crandall, Engineering Analysis, McGraw-Hill, 1956, pp. 149, 233, 372.

  Gives example to illustrate the method described herein on page 25.
- 23. Duncan, "Applications of the Galerkin method to the torsion and flexure of cylinders and prisims", Philosophical Magazine, Series 7, Vol. 25, 1938, pp. 634-649.

Applies Galerkin method to several examples of Poissons equation with polynomials as coordinate functions.

24. Finlayson and Scriven, "The method of weighted residuals and its relation to certain variational principles for the analysis of transport processes", Chemical Engineering Sciences, Vol. 20, 1965, pp. 305-404.

Discusses the method of weighted residuals and works a few examples.

25. Kaplan, "Some new methods of flux synthesis", <u>Nuclear Science and Engineering</u>, Vol. 13, 1962, pp. 22-31.

Applies Galerkin and related methods to diffusion equations subject to continuous flux and current and homogeneous boundary conditions.

26. Karnagnti, "The critical Reynold's number for the flow past a sphere", <u>Journal of the Physical Society of Japan</u>, Vol. 10, 1955, pp. 694-699.

Uses method of weighted residuals to solve nonlinear fourth order partial differential equation with nonhomogeneous boundary conditions using sines as coordinate functions and Legendre functions as weight functions.

27. Nemat-Nasser, Prasad, Henmann, "Destabalizing effect of velocity dependent forces in nonconservative continuous systems", AIAA

Journal, Vol. 4, 7, July 1966, pp. 1276-1280.

Compares an exact solution of a fourth order partial differential equation to a two term Galerkin approximation.

28. Snyder and Stewart, "Velocity and pressure profiles for newtonian creeping flow in regular packed beds of spheres", AICLE Journal, Vol. 12, 1966, pp. 167-173.

Applies Galerkin method to a homogeneous system of linear partial differential equations with nonhomogeneous boundary conditions using sines and cosines as coordinate functions.

29. Snyder, Spriggs and Stewart, "Solutions of the equations of change by Galerkins method", <u>AICLE Journal</u>, 10, 1964, pp. 535-539.

Has discussion of methods in general with no proof of his recommendations. Uses Galerkin method on nonlinear partial differential equation example with known solution to conclude results are good.

30. Toos, "The Galerkin's vector for the dynamic problems of an elastic isotropic and nonhomogeneous body", Revue Romaine de Mathematique Pures et Appliques, Vol. 10, 6, pp. 855-861.

Applies the Galerkin method to a fourth order linear nonhomogeneous partial differential equation.

#### C. Method Studies

31. Ames, "Nonlinear partial differential equations in engineering",

Mathematics in Science and Engineering, Vol. 18, Academic Press,

1965, pp. 243-270.

Discusses several variations of the method of weighted residuals. Covered herein on page 24.

32. Bellman and Richardson, "Linearization based on differential approximation and Galerkins method", Quarterly of Applied Mathematics, Vol. 24, April 1966.

Uses a one term approximation to help in gaining information to approximate a nonlinear equation in a very brief example.

- 33. Cesari, "Functional analysis and Galerkins method", Michigan

  Mathematical Journal, Vol. 11, 1964, pp. 385-414.

  See page 90 of this report.
- 34. Crandall, <u>Engineering Analysis</u>, McGraw-Hill, 1956.

  Discusses Galerkin method and applies to very simple examples.
- 35. Duncan, "The principles of the Galerkin method", TR Aeronautical Research Commission, Great Britian Air Ministry, Vol. 2, 1848, 1938, pp. 589-612.

Applies the Galerkin method to several problems quite similar to [7].

- 36. Finlayson and Scriven, "The method of weighted residuals a review",

  Applied Mechanics Reviews, Vol. 19, 9, Sept. 1966, pp. 735-747.

  Outlines the relation between Galerkins method and variational methods.
- 37. Kantorivich and Krylov, Approximate Methods of Higher Analysis, Wiley and Sons, 1964, pp. 258-304.

Gives good discussion of several approximate methods and solves several examples.

38. Mikhlin, <u>Variational Methods in Mathematical Physics</u>, Pergamon Press, 1964.

Discusses Galerkin methods and proves several convergence theorems stated herein.

39. Milne, <u>Numerical Solution of Differential Equations</u>, Wiley and Sons, 1953, pp. 114-116.

Describes the Galerkin method for nonhomogeneous boundary conditions.

40. Reiss, "Variational and related methods", bound as supplement to

Problems of Mathematical Physics by Lebeder, Prentice Hall,

1965, pp. 404-407.

Explains several variations of the Galerkin method with suggestions but offers no proof. See Chapter III.

41. Shulesko, "A new method of solving boundary value problems of mathematical physics", <u>Australian Journal of Applied Science</u>, Vol. 10, 1959, pp. 1-7.

For a discussion of his suggestion see page 23.

42. Urabe, "Galerkins procedure for nonlinear periodic systems", Archive

for Rational Mechanics and Analysis, Vol. 20, 1965, pp. 120-152.

See page 90 of this report.

- 43. Urabe and Reiter, "Numerical computation of nonlinear forced oscillations by Galerkins procedure", <u>Journal of Mathematical Analysis</u>

  and <u>Applications</u>, Vol. 14, 1966, pp. 107-140.

  See page 91 of this report.
  - 44. Urabe, "Periodic solutions of differential systems, Galerkins procedure and the method of averaging", <u>Journal of Differential</u>

    <u>Equations</u>, Vol. 2, 1966, pp. 265-280.

    See page 91 of this report.
  - 45. Babuska, <u>Numerical Processes in Differential Equations</u>, (New York: Interscience Publishers, 1966).
  - 46. Galerkin, "Rods and plates. Series occuring in various questions concerning the elastic equilibrium of rods and plates," Engineers Bulletin (Vestnik inzhenerov) 19, 1915, 897-908.
  - 47. Green, "An Expansion Method for Parabolic Partial Differential Equations," <u>Journal of Research of the National Bureau of Standards</u>, LI (September 1953), 127-132.
  - 48. Hildebrand, <u>Methods of Applied Mathematics</u>. (Englewood Cliffs: Prentice Hall, 1952), p. 286.
  - 49. Kantorovich, <u>Functional Analysis in Normed Spaces</u>, (New York: The MacMillan Co., 1964).
  - 50. Locker, "An existence analysis for nonlinear equations in Hilbert space", <u>Transactions of the American Mathematical Society</u>, Vol. 128, 3, Sept. 1967, pp. 403-413.

- 51. Milhlin, "Some sufficient conditions for the convergence of Galerkin's method," Uch. zap. Len. gos. un-ta, 135, ser. matem. nauk, 21, 1950. (restated in [38])
- 52. Petrov, "Application of Galerkin's method to the problem of stability of the flow of a viscous liquid", PMM (Applied Mathematics and Mechanics), Vol. 4, 1940, 3-12.
- 53. Petryshyn, "On the eigenvalue problem Tu Su = 0 with unbounded and nonsymmetric operators T and S," TID-4500, NYU, 1963.
- 54. Petryshyn, "Direct and iterative methods for the solution of linear operator equations in Hilbert space," Trans. Amer. Math. Soc., Vol. 105, 1962, 136-175.
- 55. Repman, "A problem in the mathematical bases of Galerkin's method for solving problems on the stability of elastic systems,"

  PMM (Applied Mathematics and Mechanics) 4, 1940.